

PCTWORLD INTELLECTUAL PROPERTY ORGANIZATION
International Bureau

INTERNATIONAL APPLICATION PUBLISHED UNDER THE PATENT COOPERATION TREATY (PCT)

(51) International Patent Classification ⁶ : C12Q 1/04, 1/68, C12N 15/10		A2	(11) International Publication Number: WO 96/17951
			(43) International Publication Date: 13 June 1996 (13.06.96)
(21) International Application Number: PCT/GB95/02875			(81) Designated States: AL, AM, AT, AU, BB, BG, BR, BY, CA, CH, CN, CZ, DE, DK, EE, ES, FI, GB, GE, HU, IS, JP, KE, KG, KP, KR, KZ, LK, LR, LS, LT, LU, LV, MD, MG, MK, MN, MW, MX, NO, NZ, PL, PT, RO, RU, SD, SE, SG, SI, SK, TJ, TM, TT, UA, UG, US, UZ, VN, European patent (AT, BE, CH, DE, DK, ES, FR, GB, GR, IE, IT, LU, MC, NL, PT, SE), OAPI patent (BF, BJ, CF, CG, CI, CM, GA, GN, ML, MR, NE, SN, TD, TG), ARIPO patent (KE, LS, MW, SD, SZ, UG).
(22) International Filing Date: 11 December 1995 (11.12.95)			
(30) Priority Data:			
9424921.6 9 December 1994 (09.12.94) GB			
9501881.8 31 January 1995 (31.01.95) GB			
9509239.1 5 May 1995 (05.05.95) GB			
(71) Applicant (for all designated States except US): RPMS TECHNOLOGY LIMITED [GB/GB]; Commonwealth Building, Du Cane Road, London W12 0NN (GB).			Published <i>Without international search report and to be republished upon receipt of that report.</i>
(72) Inventor; and			
(75) Inventor/Applicant (for US only): HOLDEN, David, William [GB/GB]; Dept. of Infectious Diseases and Bacteriology, Royal Postgraduate Medical School, Hammersmith Hospital, Du Cane Road, London W12 0NN (GB).			
(74) Agent: BASSETT, Richard; Eric Potter Clarkson, St. Mary's Court, St. Mary's Gate, Nottingham NG1 1LE (GB).			
(54) Title: IDENTIFICATION OF GENES			
(57) Abstract <p>A method for identifying a microorganism having a reduced adaptation to a particular environment comprising the steps of: (1) providing a plurality of microorganisms each of which is independently mutated by the insertional inactivation of a gene with a nucleic acid comprising a unique marker sequence so that each mutant contains a different marker sequence, or clones of the said microorganism; (2) providing individually a stored sample of each mutant produced by step (1) and providing individually stored nucleic acid comprising the unique marker sequence from each individual mutant; (3) introducing a plurality of mutants produced by step (1) into the said particular environment and allowing those microorganisms which are able to do so to grow in the said environment; (4) retrieving microorganisms from the said environment or a selected part thereof and isolating the nucleic acid from the retrieved microorganisms; (5) comparing any marker sequences in the nucleic acid isolated in step (4) to the unique marker sequence of each individual mutant stored as in step (2); and (6) selecting an individual mutant which does not contain any of the marker sequences as isolated in step (4).</p>			

FOR THE PURPOSES OF INFORMATION ONLY

Codes used to identify States party to the PCT on the front pages of pamphlets publishing international applications under the PCT.

AT	Austria	GB	United Kingdom	MR	Mauritania
AU	Australia	GE	Georgia	MW	Malawi
BB	Barbados	GN	Guinea	NE	Niger
BE	Belgium	GR	Greece	NL	Netherlands
BF	Burkina Faso	HU	Hungary	NO	Norway
BG	Bulgaria	IE	Ireland	NZ	New Zealand
BJ	Benin	IT	Italy	PL	Poland
BR	Brazil	JP	Japan	PT	Portugal
BY	Belarus	KE	Kenya	RO	Romania
CA	Canada	KG	Kyrgyzstan	RU	Russian Federation
CF	Central African Republic	KP	Democratic People's Republic of Korea	SD	Sudan
CG	Congo	KR	Republic of Korea	SE	Sweden
CH	Switzerland	KZ	Kazakhstan	SI	Slovenia
CI	Côte d'Ivoire	LI	Liechtenstein	SK	Slovakia
CM	Cameroon	LK	Sri Lanka	SN	Senegal
CN	China	LU	Luxembourg	TD	Chad
CS	Czechoslovakia	LV	Latvia	TG	Togo
CZ	Czech Republic	MC	Monaco	TJ	Tajikistan
DE	Germany	MD	Republic of Moldova	TT	Trinidad and Tobago
DK	Denmark	MG	Madagascar	UA	Ukraine
ES	Spain	ML	Mali	US	United States of America
FI	Finland	MN	Mongolia	UZ	Uzbekistan
FR	France			VN	Viet Nam
GA	Gabon				

IDENTIFICATION OF GENES

The present invention relates to methods for the identification of genes involved in the adaptation of a microorganism to its environment,
5 particularly the identification of genes responsible for the virulence of a pathogenic microorganism.

Background to the invention

10 Antibiotic resistance in bacterial and other pathogens is becoming increasingly important. It is therefore important to find new therapeutic approaches to attack pathogenic microorganisms.

Pathogenic microorganisms have to evade the host's defence mechanisms
15 and be able to grow in a poor nutritional environment to establish an infection. To do so a number of "virulence" genes of the microorganism are required.

Virulence genes have been detected using classical genetics and a variety
20 of approaches have been used to exploit transposon mutagenesis for the identification of bacterial virulence genes. For example, mutants have been screened for defined physiological defects, such as the loss of iron regulated proteins (Holland *et al*, 1992), or in assays to study the penetration of epithelial cells (Finlay *et al*, 1988) and survival within
25 macrophages (Fields *et al*, 1989; Miller *et al*, 1989a; Groisman *et al*, 1989). Transposon mutants have also been tested for altered virulence in live animal models of infection (Miller *et al*, 1989b). This approach has the advantage that genes can be identified which are important during different stages of infection, but is severely limited by the need to test a
30 wide range of mutants individually for alterations to virulence. Miller *et*

al (1989b) used groups of 8 to 10 mice and infected orally 95 separate groups with a different mutant thereby using between 760 and 950 mice. Because of the extremely large numbers of animals required, comprehensive screening of a bacterial genome for virulence genes has not
5 been feasible.

Recently a genetic system (*in vivo* expression technology [IVET]) was described which positively selects for *Salmonella* genes that are specifically induced during infection (Mahan *et al*, 1993). The technique
10 will identify genes that are expressed at a particular stage in the infection process. However, it will not identify virulence genes that are regulated posttranscriptionally, and more importantly, will not provide information on whether the gene(s) which have been identified are actually required for, or contribute to, the infection process.

15

Lee & Falkow (1994) *Methods Enzymol.* 236, 531-545 describe a method of identifying factors influencing the invasion of *Salmonella* into mammalian cells *in vitro* by isolating hyperinvasive mutants.

20 Walsh and Cepko (1992) *Science* 255, 434-440 describe a method of tracking the spatial location of cerebral cortical progenitor cells during the development of the cerebral cortex in the rat. The Walsh and Cepko method uses a tag that contains a unique nucleic acid sequence and the lacZ gene but there is no indication that useful mutants or genes could be
25 detected by their method.

WO 94/26933 and Smith *et al* (1995) *Proc. Natl. Acad. Sci. USA* 92, 6479-6483 describe methods aimed at the identification of the functional regions of a known gene, or at least of a DNA molecule for which some
30 sequence information is available.

Groisman *et al* (1993) *Proc. Natl. Acad. Sci. USA* 90, 1033-1037 describes the molecular, functional and evolutionary analysis of sequences specific to *Salmonella*.

5 Some virulence genes are already known for pathogenic microorganisms such as *Escherichia coli*, *Salmonella typhimurium*, *Salmonella typhi*, *Vibrio cholerae*, *Clostridium botulinum*, *Yersinia pestis*, *Shigella flexneri* and *Listeria monocytogenes* but in all cases only a relatively small number of the total have been identified.

10

The disease which *Salmonella typhimurium* causes in mice provides a good experimental model of typhoid fever (Carter & Collins, 1974). Approximately forty two genes affecting *Salmonella* virulence have been identified to date (Groisman & Ochman, 1994). These represent
15 approximately one third of the total number of predicted virulence genes (Groisman and Saier, 1990).

The object of the present invention is to identify genes involved in the adaptation of a microorganism to its environment, particularly to identify
20 further virulence genes in pathogenic microorganisms, with increased efficiency. A further object is to reduce the number of experimental animals used in identifying virulence genes. Still further objects of the invention provide vaccines, and methods for screening for drugs which reduce virulence.

25

Summary of the invention

A first aspect of the invention provides a method for identifying a microorganism having a reduced adaptation to a particular environment
30 comprising the steps of:

- (1) providing a plurality of microorganisms each of which is independently mutated by the insertional inactivation of a gene with a nucleic acid comprising a unique marker sequence so that each mutant contains a different marker sequence, or clones of the said microorganism;
- 5 (2) providing individually a stored sample of each mutant produced by step (1) and providing individually stored nucleic acid comprising the unique marker sequence from each individual mutant;
- (3) introducing a plurality of mutants produced by step (1) into the said particular environment and allowing those microorganisms which
- 10 are able to do so to grow in the said environment;
- (4) retrieving microorganisms from the said environment or a selected part thereof and isolating the nucleic acid from the retrieved microorganisms;
- (5) comparing any marker sequences in the nucleic acid isolated
- 15 in step (4) to the unique marker sequence of each individual mutant stored as in step (2); and
- (6) selecting an individual mutant which does not contain any of the marker sequences as isolated in step (4).
- 20 Thus, the method uses negative selection to identify microorganisms with reduced capacity to proliferate in the environment.

A microorganism can live in a number of different environments and it is known that particular genes and their products allow the microorganism

25 to adapt to a particular environment. For example, in order for a pathogenic microorganism, such as a pathogenic bacterium or pathogenic fungus, to survive in its host the product of one or more virulence genes is required. Thus, in a preferred embodiment of the invention a gene of a microorganism which allows the microorganism to adapt to a particular

30 environment is a virulence gene.

Conveniently, the particular environment is a differentiated multicellular organism such as a plant or animal. Many bacteria and fungi are known to infect plants and they are able to survive within the plant and cause disease because of the presence of and expression from virulence genes.

- 5 Suitable microorganisms when the particular environment is a plant include the bacteria *Agrobacterium tumefaciens* which forms tumours (galls) particularly in grape; *Erwinia amylovora*; *Pseudomonas solanacearum* which causes wilt in a wide range of plants; *Rhizobium leguminosarum* which causes disease in beans; *Xanthomonas campestris*
10 p.v. *citri* which causes canker in citrus fruits; and include the fungi *Magnaporthe grisea* which causes rice blast disease; *Fusarium* spp. which cause a variety of plant diseases; *Erysiphe* spp.; *Colletotrichum gloeosporioides*; *Gaeumannomyces graminis* which causes root and crown diseases in cereals and grasses; *Glomus* spp., *Laccaria* spp.; *Leptosphaeria*
15 *maculans*; *Phoma tracheiphila*; *Phytophthora* spp., *Pyrenophora teres*; *Verticillium alboatrum* and *V. dahliae*; and *Mycosphaerella musicola* and *M. fijiensis*. As described in more detail below, when the microorganism is a fungus a haploid phase to its life cycle is required.

- 20 Similarly, many microorganisms including bacteria, fungi, protozoa and trypanosomes are known to infect animals, particularly mammals including humans. Survival of the microorganism within the animal and the ability of the microorganism to cause disease is due in large part to the presence of and expression from virulence genes. Suitable bacteria include
25 *Bordetella* spp. particularly *B. pertussis*, *Campylobacter* spp. particularly *C. jejuni*, *Clostridium* spp. particularly *C. botulinum*, *Enterococcus* spp. particularly *E. faecalis*, *Escherichia* spp. particularly *E. coli*, *Haemophilus* spp. particularly *H. ducreyi* and *H. influenzae*, *Helicobacter* spp. particularly *H. pylori*, *Klebsiella* spp. particularly *K. pneumoniae*,
30 *Legionella* spp. particularly *L. pneumophila*, *Listeria* spp. particularly *L.*

monocytogenes, *Mycobacterium* spp. particularly *M. smegmatis* and *M. tuberculosis*, *Neisseria* spp. particularly *N. gonorrhoeae* and *N. meningitidis*, *Pseudomonas* spp., particularly *Ps. aeruginosa*, *Salmonella* spp., *Shigella* spp., *Staphylococcus* spp. particularly *S. aureus*,
5 *Streptococcus* spp. particularly *S. pyogenes* and *pneumoniae*, *Vibrio* spp. and *Yersinia* spp. particularly *Y. pestis*. All of these bacteria cause disease in man and also there are animal models of the disease. Thus, when these bacteria are used in the method of the invention, the particular environment is an animal which they can infect and in which they cause
10 disease. For example, when *Salmonella typhimurium* is used to infect a mouse the mouse develops a disease which serves as a model for typhoid fever in man. *Staphylococcus aureus* causes bacteraemia and renal abscess formation in mice (Albus *et al* (1991) *Infect. Immun.* 59, 1008-1014) and endocarditis in rabbits (Perlman & Freedman (1971) *Yale J. Biol. Med.*
15 44, 206-213).

It is required that a fungus or higher eukaryotic parasite is haploid for the relevant parts of its life (such as growth in the environment). Preferably, a DNA-mediated integrative transformation system is available and, when
20 the microorganism is a human pathogen, conveniently an animal model of the human disease is available. Suitable fungi pathogenic to humans include certain *Aspergillus* spp. (for example *A. fumigatus*), *Cryptococcus neoformans* and *Histoplasma capsulatum*. Clearly the above-mentioned fungi have a haploid phase and a DNA-mediated integrative transformation
25 systems are available for them. *Toxoplasma* may also be used, being a parasite with a haploid phase during infection. Bacteria have a haploid genome.

Animal models of human disease are often available in which the animal
30 is a mouse, rat, rabbit, dog or monkey. It is preferred if the animal is a

mouse. Virulence genes detected by the method of the invention using an animal model of a human disease are clearly very likely to be genes which determine the virulence of the microorganism in man.

- 5 Particularly preferred microorganisms for use in the methods of the invention are *Salmonella typhimurium*, *Staphylococcus aureus*, *Streptococcus pneumoniae*, *Enterococcus faecalis*, *Pseudomonas aeruginosa* and *Aspergillus fumigatus*.

- 10 A preferred embodiment of the invention is now described.

A nucleic acid comprising a unique marker sequence is made as follows. A complex pool of double stranded DNA sequence "tags" is generated using oligonucleotide synthesis and a polymerase chain reaction (PCR).

- 15 Each DNA "tag" has a unique sequence of between about 20 and 80 bp, preferably about 40 bp which is flanked by "arms" of about 15 to 30 bp, preferably about 20 bp, which are common to all "tags". The number of bp in the unique sequence is sufficient to allow large numbers (for example $> 10^{10}$) of unique sequences to be generated by random
20 oligonucleotide synthesis but not too large to allow the formation of secondary structures which may interfere with a PCR. Similarly, the length of the arms should be sufficient to allow efficient priming of oligonucleotides in a PCR.

- 25 It is well known that the sequence at the 5' end of the oligonucleotide need not match the target sequence to be amplified.

- It is usual that the PCR primers do not contain any complementary structures with each other longer than 2 bases, especially at their 3' ends,
30 as this feature may promote the formation of an artifactual product called

“primer dimer”. When the 3' ends of the two primers hybridize, they form a “primed template” complex, and primer extension results in a short duplex product called “primer dimer”.

- 5 Internal secondary structure should be avoided in primers. For symmetric PCR, a 40-60% G+C content is often recommended for both primers, with no long stretches of any one base. The classical melting temperature calculations used in conjunction with DNA probe hybridization studies often predict that a given primer should anneal at a specific temperature
10 or that the 72°C extension temperature will dissociate the primer/template hybrid prematurely. In practice, the hybrids are more effective in the PCR process than generally predicted by simple T_m calculations.

Optimum annealing temperatures may be determined empirically and may
15 be higher than predicted. *Taq* DNA polymerase does have activity in the 37-55°C region, so primer extension will occur during the annealing step and the hybrid will be stabilized. The concentrations of the primers are equal in conventional (symmetric) PCR and, typically, within 0.1- to 1- μ M range.

- 20 The “tags” are ligated into a transposon or transposon-like element to form the nucleic acid comprising a unique marker sequence. Conveniently, the transposon is carried on a suicide vector which is maintained as a plasmid in a “helper” organism, but which is lost after
25 transfer to the microorganism of the method of the invention. For example, the “helper” organism may be a strain of *Escherichia coli*, the microorganism of the method may be *Salmonella* and the transfer is a conjugal transfer. Although the transposon can be lost after transfer, in a proportion of cells it undergoes a transposition event through which it
30 integrates at random, along with its unique tag, into the genome of the

microorganism used in the method. It is most preferred if the transposon or transposon-like element can be selected. For example, in the case of *Salmonella*, a kanamycin resistance gene may be present in the transposon and exconjugants are selected on medium containing kanamycin. It is also possible to complement an auxotrophic marker in the recipient cell with a functional gene in the nucleic acid comprising the unique marker. This method is particularly convenient when fungi are used in the method.

Preferably the complementing functional gene is not derived from the same species as the recipient microorganism, otherwise non-random integration events may occur.

It is also particularly convenient if the transposon or transposon-like element is carried on a vector which is maintained episomally (ie not as part of the chromosome) in the microorganism used in the method of the first aspect of the invention when in a first given condition whereas, upon changing that condition to a second given condition, the episome cannot be maintained permitting the selection of a cell in which the transposon or transposon-like element has undergone a transposition event through which it integrates at random, along with its unique tag, into the genome of the microorganism used in the method. This particularly convenient embodiment is advantageous because, once a microorganism carrying the episomal vector is made, then each time the transposition event is selected for or induced by changing the condition of the microorganism (or a clone thereof) from a first given condition to a second given condition, the transposon can integrate at a different site in the genome of the microorganism. Thus, once a master collection of microorganisms are made, each member of which contains a unique tag sequence in the transposon or transposon-like element carried on the episomal vector (when in the first given condition), it can be used repeatedly to generate

pools of random insertional mutants, each of which contains a different tag sequence (ie unique within the pool). This embodiment is particularly useful because (a) it reduces the number and complexity of manipulations required to generate the plurality ("pool") of independently mutated
5 microorganisms in step (1) of the method; and (b) the number of different tags need only be the same as the number of microorganisms in the plurality of microorganisms in step (1) of the method. Point (a) makes the method easier to use in organisms for which transposon mutagenesis is more difficult to perform (for example, *Staphylococcus aureus*) and point
10 (b) means that tag sequences with particularly good hybridisation characteristics can be selected therefore making quality control easier. As is described in more detail below, the "pool" size is conveniently about 100 or 200 independently-mutated microorganisms and, therefore the master collection of microorganisms is conveniently stored in one or two
15 96-well microtitre plates.

In a particularly preferred embodiment the first given condition is a first particular temperature or temperature range such as 25°C to 32°C, most preferably about 30°C and the second given condition is a second
20 particular temperature or temperature range such as 35°C to 45°C, most preferably 42°C. In further preferred embodiments the first given condition is the presence of an antibiotic, such as streptomycin, and the second given condition is the absence of the said antibiotic; or the first given condition is the absence of an antibiotic and the second given
25 condition is the presence of the said antibiotic.

Transposons suitable for integration into the genome of Gram negative bacteria include Tn5, Tn10 and derivatives thereof. Transposons suitable for integration into the genome of Gram positive bacteria include Tn916
30 and derivatives or analogues thereof. Transposons particularly suited for

use with *Staphylococcus aureus* include Tn917 (Cheung *et al* (1992) *Proc. Natl. Acad. Sci. USA* 89, 6462-6466) and Tn918 (Albus *et al* (1991) *Infect. Immun.* 59, 1008-1014).

- 5 It is particularly preferred if the transposons have the properties of the Tn917 derivatives described by Camilli *et al* (1990) *J. Bacteriol.* 172, 3738-3744, and are carried by a temperature-sensitive vector such as pE194Ts (Villafane *et al* (1987) *J. Bacteriol.* 169, 4822-4829).
- 10 It will be appreciated that although transposons are convenient for insertionally inactivating a gene, any other known method, or method developed in the future may be used. A further convenient method of insertionally inactivating a gene, particularly in certain bacteria such as *Streptococcus*, is using insertion-duplication mutagenesis such as that
- 15 described in Morrison *et al* (1984) *J. Bacteriol.* 159, 870 with respect to *S. pneumoniae*. The general method may also be applied to other microorganisms, especially bacteria.

For fungi, insertional mutations are created by transformation using DNA

20 fragments or plasmids carrying the "tags" and, preferably, selectable markers encoding, for example, resistance to hygromycin B or phleomycin (see Smith *et al* (1994) *Infect. Immunol.* 62, 5247-5254). Random, single integration of DNA fragments encoding hygromycin B resistance into the genome of filamentous fungi, using restriction enzyme mediated

25 integration (REMI; Schiestl & Petes (1991); Lu *et al* (1994) *Proc. Natl. Acad. Sci. USA* 91, 12649-12653) are known.

A simple insertional mutagenesis technique for a fungus is described in Schiestl & Petes (1994) incorporated herein by reference, and include, for

30 example, the use of Ty elements and ribosomal DNA in yeast.

Random integration of the transposon or other DNA sequence allows isolation of a plurality of independently mutated microorganisms wherein a different gene is insertionally inactivated in each mutant and each mutant contains a different marker sequence.

5

A library of such insertion mutants is arrayed in well microtitre dishes so that each well contains a different mutant microorganism. DNA comprising the unique marker sequence from each individual mutant microorganism (conveniently, the total DNA from the clone is used) is stored. Conveniently, this is done by removing a sample of the
10 microorganism from the microtitre dish, spotting it onto a nucleic acid hybridisation membrane (such as nitrocellulose or nylon membranes), lysing the microorganism in alkali and fixing the nucleic acid to the membrane. Thus, a replica of the contents of the well microtitre dishes
15 is made.

Pools of the microorganisms from the well microtitre dish are made and DNA is extracted. This DNA is used as a target for a PCR using primers that anneal to the common "arms" flanking the "tags" and the amplified
20 DNA is labelled, for example with ^{32}P . The product of the PCR is used to probe the DNA stored from each individual mutant to provide a reference hybridisation pattern for the replicas of the well microtitre dishes. This is a check that each of the individual microorganisms does, in fact, contain a marker sequence and that the marker sequence can be
25 amplified and labelled efficiently.

Pools of transposon mutants are made to introduce into the particular environment. Conveniently, 96-well microtitre dishes are used and the pool contains 96 transposon mutants. However, the lower limit for the
30 pool is two mutants: there is no theoretical upper limit to the size of the

pool but, as discussed below, the upper limit may be determined in relation to the environment into which the mutants are introduced.

Once the microorganisms are introduced into the said particular environment those microorganisms which are able to do so are allowed to grow in the environment. The length of time the microorganisms are left in the environment is determined by the nature of the microorganism and the environment. After a suitable length of time, the microorganisms are recovered from the environment, DNA is extracted and the DNA is used as a template for a PCR using primers that anneal to the "arms" flanking the "tags". The PCR product is labelled, for example with ^{32}P , and is used to probe the DNA stored from each individual mutant replicated from the wellled microtitre dish. Stored DNA are identified which hybridise weakly or not at all with the probe generated from the DNA isolated from the microorganisms recovered from environment. These non-hybridising DNAs correspond to mutants whose adaptation to the particular environment has been attenuated by insertion of the transposon or other DNA sequence.

In a particularly preferred embodiment the "arms" have no, or very little, label compared to the "tags". For example, the PCR primers are suitably designed to contain no, or a single, G residue, the ^{32}P -labelled nucleotide is dCTP and, in this case, no or one radiolabelled C residue is incorporated in each "arm" but a greater number of radiolabelled C residues are incorporated in the "tag". It is preferred if the "tag" has at least ten-fold more label incorporated than the "arms"; preferably twenty-fold or more; more preferably fifty-fold or more. Conveniently the "arms" can be removed from the "tag" using a suitable restriction enzyme, a site for which may be incorporated in the primer design.

As discussed above, a particularly preferred embodiment of the invention is when the microorganism is a pathogenic microorganism and the particular environment is an animal. In this embodiment, the size of the pool of mutants introduced into the animal is determined by (a) the number of cells of each mutant that are likely to survive in the animal (assuming a virulence gene has not been inactivated) and (b) the total inoculum of the microorganism. If the number in (a) is too low then false positive results may arise and if the number in (b) is too high then the animal may die before enough mutants have had a chance to grow in the desired way. The number of cells in (a) can be determined for each microorganism used but it is preferably more than 50, more preferably more than 100.

The number of different mutants that can be introduced into a single animal is preferably between 50 and 500, conveniently about 100. It is preferred if the total inoculum does not exceed 10^6 cells (and it is preferably 10^5 cells) although the size of the inoculum may be varied above or below this amount depending on the microorganism and the animal.

In a particularly convenient method an inoculum of 10^5 is used containing 1000 cells each of 100 different mutants for a single animal. It will be appreciated that in this method one animal can be used to screen 100 mutants compared to prior art methods which require at least 100 animals to screen 100 mutants.

However, it is convenient to inoculate three animals with the same pool of mutants so that at least two can be investigated (one as a replica to check the reliability of the method), whilst the third is kept as a back-up. Nevertheless, the method still provides a greater than 30-fold saving in the

number of animals used.

The time between the pool of mutants being introduced into the animal and the microorganisms being recovered may vary with the microorganism and animal used. For example, when the animal is a mouse and the
5 microorganism is *Salmonella typhimurium*, the time between inoculation and recovery is about three days.

In one embodiment of the invention microorganisms are retrieved from the
10 environment in step (5) at a site remote from the site of introduction in step (4), so that the virulence genes being investigated include those involved in the spread of the microorganism between the two sites.

For example, in a plant the microorganism may be introduced in a lesion
15 in the stem or at one site on a leaf and the microorganism retrieved from another site on the leaf where a disease state is indicated.

In the case of an animal, the microorganism may be introduced orally, intraperitoneally, intravenously or intranasally and is retrieved at a later
20 time from an internal organ such as the spleen. It may be useful to compare the virulence genes identified by oral administration and those identified by intraperitoneal administration as some genes may be required to establish infection by one route but not by the other. It is preferred if *Salmonella* is introduced intraperitoneally.

25

Other preferred environments which may be used to identify virulence genes are animal cells in culture (particularly macrophages and epithelial cells) and plant cells in culture. Although using cells in culture will be useful in its own right, it will also complement the use of the whole
30 animal or plant, as the case may be, as the environment.

It is also preferred if the environment is a part of the animal body. Within a given host-parasite interaction, a number of different environments are possible, including different organs and tissues, and parts thereof such as the Peyer's patch.

5

The number of individual microorganisms (ie cells) recovered from the environment should be at least twice, preferably at least ten times, more preferably 100-times the number of different mutants introduced into the environment. For example, when an animal is inoculated with 100
10 different mutants around 10,000 individual microorganisms should be retrieved and their marker DNA isolated.

A further preferred embodiment comprises the steps:

15 (1A) removing auxotrophs from the plurality of mutants produced in step (1); or

(6A) determining whether the mutant selected in step (6) is an auxotroph;
or

20

both (1A) and (6A).

It is desirable to distinguish an auxotroph (that is a mutant microorganism which requires growth factors not needed by the wild type or by
25 prototrophs) and a mutant microorganism wherein a gene allowing the microorganism to adapt to a particular environment is inactivated. Conveniently, this is done between steps (1) and (2) or after step (6).

Preferably auxotrophs are not removed when virulence genes are being
30 identified.

A second aspect of the invention provides a method of identifying a gene which allows a microorganism to adapt to a particular environment, the method comprising the method of the first aspect of the invention, followed by the additional step:

5

(7) isolating the insertionally-inactivated gene or part thereof from the individual mutant selected in step (6).

Methods for isolating a gene containing a unique marker are well known
10 in the art of molecular biology.

A further preferred embodiment comprises the further additional step:

(8) isolating from a wild-type microorganism the corresponding wild-
15 type gene using the insertionally-inactivated gene isolated in step (7) or part thereof as a probe.

Methods for gene probing are well known in the art of molecular biology.

20 Molecular biological methods suitable for use in the practice of the present invention are disclosed in Sambrook *et al* (1989) incorporated herein by reference.

When the microorganism is a microorganism pathogenic to an animal and
25 the gene is a virulence gene and a transposon has been used to insertionally inactivate the gene, it is convenient for the virulence genes to be cloned by digesting genomic DNA from the individual mutant selected in step (6) with a restriction enzyme which cuts outside the transposon, ligating size-fractionated DNA containing the transposon into
30 a plasmid, and selecting plasmid recombinants on the basis of antibiotic

resistance conferred by the transposon and not by the plasmid. The microorganism genomic DNA adjacent to the transposon is sequenced using two primers which anneal to the terminal regions of the transposon, and two primers which anneal close to the polylinker sequences of the plasmid. The sequences may be subjected to DNA database searches to determine if the transposon has interrupted a known virulence gene. Thus, conveniently, sequence obtained by this method is compared against the sequences present in the publicly available databases such as EMBL and GenBank. Finally, if the interrupted sequence appears to be in a new virulence gene, the mutation is transferred to a new genetic background (for example by phage P22-mediated transduction in the case of *Salmonella*) and the LD₅₀ of the mutant strain is determined to confirm that the avirulent phenotype is due to the transposition event and not a secondary mutation.

15

The number of individual mutants screened in order to detect all of the virulence genes in a microorganism depends on the number of genes in the genome of the microorganism. For example, it is likely that 3000-5000 mutants of *Salmonella typhimurium* need to be screened in order to detect the majority of virulence genes whereas for *Aspergillus* spp., which has a larger genome than *Salmonella*, around 20 000 mutants are screened. Approximately 4% of non-essential *S. typhimurium* genes are thought to be required for virulence (Grossman & Saier, 1990) and, if so, the *S. typhimurium* genome contains approximately 150 virulence genes. However, the methods of the invention provide a faster, more convenient and much more practicable route to identifying virulence genes.

25

A third aspect of the invention provides a microorganism obtained using the method of the first aspect of the invention.

30

Such microorganisms are useful because they have the property of not being adapted to survive in a particular environment.

In a preferred embodiment, a pathogenic microorganism is not adapted to survive in a host organism (environment) and, in the case of microorganisms that are pathogenic to animals, particularly mammals, more particularly humans, the mutant obtained by the method of the invention may be used in a vaccine. The mutant is avirulent, and therefore expected to be suitable for administration to a patient, but it is expected to be antigenic and give rise to a protective immune response.

In a further preferred embodiment the pathogenic microorganism not adapted to survive in a host organism, obtained by the methods of the invention, is modified, preferably by the introduction of a suitable DNA sequence, to express an antigenic epitope from another pathogen. This modified microorganism can act as a vaccine for that other pathogen.

A fourth aspect of the invention provides a microorganism comprising a mutation in a gene identified using the method of the second aspect of the invention.

Thus, although the microorganism of the third aspect of the invention is useful, it is preferred if a mutation is specifically introduced into the identified gene. In a preferred embodiment, particularly when the microorganism is to be used in a vaccine, the mutation in the gene is a deletion or a frameshift mutation or any other mutation which is substantially incapable of reverting. Such gene-specific mutations can be made using standard procedures such as introducing into the microorganism a copy of the mutant gene on an autonomous replicon (such as a plasmid or viral genome) and relying on homologous

recombination to introduce the mutation into the copy of the gene in the genome of the microorganism.

Fifth and sixth aspects of the invention provide a suitable microorganism
5 for use in a vaccine and a vaccine comprising a suitable microorganism
and a pharmaceutically-acceptable carrier.

The suitable microorganism is the aforementioned avirulent mutant.

10 Active immunisation of the patient is preferred. In this approach, one or
more mutant microorganisms are prepared in an immunogenic formulation
containing suitable adjuvants and carriers and administered to the patient
in known ways. Suitable adjuvants include Freund's complete or
incomplete adjuvant, muramyl dipeptide, the "Iscoms" of EP 109 942, EP
15 180 564 and EP 231 039, aluminium hydroxide, saponin, DEAE-dextran,
neutral oils (such as miglyol), vegetable oils (such as arachis oil),
liposomes, Pluronic polyols or the Ribi adjuvant system (see, for example
GB-A-2 189 141). "Pluronic" is a Registered Trade Mark. The patient
to be immunised is a patient requiring to be protected from the disease
20 caused by the virulent form of the microorganism.

The aforementioned avirulent microorganisms of the invention or a
formulation thereof may be administered by any conventional method
including oral and parenteral (eg subcutaneous or intramuscular) injection.
25 The treatment may consist of a single dose or a plurality of doses over a
period of time.

Whilst it is possible for an avirulent microorganism of the invention to be
administered alone, it is preferable to present it as a pharmaceutical
30 formulation, together with one or more acceptable carriers. The carrier(s)

must be "acceptable" in the sense of being compatible with the avirulent microorganism of the invention and not deleterious to the recipients thereof. Typically, the carriers will be water or saline which will be sterile and pyrogen free.

5

It will be appreciated that the vaccine of the invention, depending on its microorganism component, may be useful in the fields of human medicine and veterinary medicine.

- 10 Diseases caused by microorganisms are known in many animals, such as domestic animals. The vaccines of the invention, when containing an appropriate avirulent microorganism, particularly avirulent bacterium, are useful in man but also in, for example, cows, sheep, pigs, horses, dogs and cats, and in poultry such as chickens, turkeys, ducks and geese.

15

Seventh and eighth aspects of the invention provide a gene obtained by the method of the second aspect of the invention, and a polypeptide encoded thereby.

- 20 By "gene" we include not only the regions of DNA that code for a polypeptide but also regulatory regions of DNA such as regions of DNA that regulate transcription, translation and, for some microorganisms, splicing of RNA. Thus, the gene includes promoters, transcription terminators, ribosome-binding sequences and for some organisms introns
- 25 and splice recognition sites.

- Typically, sequence information of the inactivated gene obtained in step 7 is derived. Conveniently, sequences close to the ends of the transposon are used as the hybridisation site of a sequencing primer. The derived
- 30 sequence or a DNA restriction fragment adjacent to the inactivated gene

itself is used to make a hybridisation probe with which to identify, and isolate from a wild-type organism, the corresponding wild type gene.

It is preferred if the hybridisation probing is done under stringent
5 conditions to ensure that the gene, and not a relative, is obtained. By
"stringent" we mean that the gene hybridises to the probe when the gene
is immobilised on a membrane and the probe (which, in this case is > 200
nucleotides in length) is in solution and the immobilised gene/hybridised
probe is washed in 0.1 x SSC at 65°C for 10 min. SSC is 0.15 M
10 NaCl/0.015 M Na citrate.

Preferred probe sequences for cloning *Salmonella* virulence genes are shown in Figures 5 and 6 and described in Example 2.

15 In a particularly preferred embodiment the *Salmonella* virulence genes comprise the sequence shown in Figures 5 and 6 and described in Example 2.

In further preference the genes are those contained within, or at least part
20 of which is contained within, the sequences shown in Figures 11 and 12
and which have been identified by the method of the second aspect of the
invention. The sequences shown in Figures 11 and 12 are part of a gene
cluster from *Salmonella typhimurium* which I have called virulence gene
cluster 2 (VGC2). The position of transposon insertions are indicated
25 within the sequence, and these transposon insertions inactivate a virulence
determinant of the organism. As is discussed more fully below, and in
particular in Example 4, when the method of the second aspect of the
invention is used to identify virulence genes in *Salmonella typhimurium*,
many of the nucleic acid insertions (and therefore genes identified) are
30 clustered in a relatively small part of the genome. This region. VGC2.

contains other virulence genes which, as described below, form part of the invention.

The gene isolated by the method of the invention can be expressed in a suitable host cell. Thus, the gene (DNA) may be used in accordance with known techniques, appropriately modified in view of the teachings contained herein, to construct an expression vector, which is then used to transform an appropriate host cell for the expression and production of the polypeptide of the invention. Such techniques include those disclosed in US Patent Nos. 4,440,859 issued 3 April 1984 to Rutter *et al*, 4,530,901 issued 23 July 1985 to Weissman, 4,582,800 issued 15 April 1986 to Crowl, 4,677,063 issued 30 June 1987 to Mark *et al*, 4,678,751 issued 7 July 1987 to Goeddel, 4,704,362 issued 3 November 1987 to Itakura *et al*, 4,710,463 issued 1 December 1987 to Murray, 4,757,006 issued 12 July 1988 to Toole, Jr. *et al*, 4,766,075 issued 23 August 1988 to Goeddel *et al* and 4,810,648 issued 7 March 1989 to Stalker, all of which are incorporated herein by reference.

The DNA encoding the polypeptide constituting the compound of the invention may be joined to a wide variety of other DNA sequences for introduction into an appropriate host. The companion DNA will depend upon the nature of the host, the manner of the introduction of the DNA into the host, and whether episomal maintenance or integration is desired.

Generally, the DNA is inserted into an expression vector, such as a plasmid, in proper orientation and correct reading frame for expression. If necessary, the DNA may be linked to the appropriate transcriptional and translational regulatory control nucleotide sequences recognised by the desired host, although such controls are generally available in the expression vector. The vector is then introduced into the host through

standard techniques. Generally, not all of the hosts will be transformed by the vector. Therefore, it will be necessary to select for transformed host cells. One selection technique involves incorporating into the expression vector a DNA sequence, with any necessary control elements, that codes for a selectable trait in the transformed cell, such as antibiotic resistance. Alternatively, the gene for such selectable trait can be on another vector, which is used to co-transform the desired host cell.

Host cells that have been transformed by the recombinant DNA of the invention are then cultured for a sufficient time and under appropriate conditions known to those skilled in the art in view of the teachings disclosed herein to permit the expression of the polypeptide, which can then be recovered.

Many expression systems are known, including bacteria (for example *E. coli* and *Bacillus subtilis*), yeasts (for example *Saccharomyces cerevisiae*), filamentous fungi (for example *Aspergillus*), plant cells, animal cells and insect cells.

The vectors include a prokaryotic replicon, such as the ColE1 *ori*, for propagation in a prokaryote, even if the vector is to be used for expression in other, non-prokaryotic, cell types. The vectors can also include an appropriate promoter such as a prokaryotic promoter capable of directing the expression (transcription and translation) of the genes in a bacterial host cell, such as *E. coli*, transformed therewith.

A promoter is an expression control element formed by a DNA sequence that permits binding of RNA polymerase and transcription to occur. Promoter sequences compatible with exemplary bacterial hosts are typically provided in plasmid vectors containing convenient restriction sites

for insertion of a DNA segment of the present invention.

Typical prokaryotic vector plasmids are pUC18, pUC19, pBR322 and pBR329 available from Biorad Laboratories, (Richmond, CA, USA) and
5 pTrc99A and pKK223-3 available from Pharmacia, Piscataway, NJ, USA.

A typical mammalian cell vector plasmid is pSVL available from Pharmacia, Piscataway, NJ, USA. This vector uses the SV40 late promoter to drive expression of cloned genes, the highest level of
10 expression being found in T antigen-producing cells, such as COS-1 cells.

An example of an inducible mammalian expression vector is pMSG, also available from Pharmacia. This vector uses the glucocorticoid-inducible promoter of the mouse mammary tumour virus long terminal repeat to
15 drive expression of the cloned gene.

Useful yeast plasmid vectors are pRS403-406 and pRS413-416 and are generally available from Stratagene Cloning Systems, La Jolla, CA 92037, USA. Plasmids pRS403, pRS404, pRS405 and pRS406 are Yeast
20 Integrating plasmids (YIps) and incorporate the yeast selectable markers *HIS3*, *TRP1*, *LEU2* and *URA3*. Plasmids pRS413-416 are Yeast Centromere plasmids (YCps)

A variety of methods have been developed to operably link DNA to
25 vectors via complementary cohesive termini. For instance, complementary homopolymer tracts can be added to the DNA segment to be inserted to the vector DNA. The vector and DNA segment are then joined by hydrogen bonding between the complementary homopolymeric tails to form recombinant DNA molecules.

Synthetic linkers containing one or more restriction sites provide an alternative method of joining the DNA segment to vectors. The DNA segment, generated by endonuclease restriction digestion as described earlier, is treated with bacteriophage T4 DNA polymerase or *E. coli* DNA
5 polymerase I, enzymes that remove protruding, 3'-single-stranded termini with their 3'-5'-exonucleolytic activities, and fill in recessed 3'-ends with their polymerizing activities.

The combination of these activities therefore generates blunt-ended DNA
10 segments. The blunt-ended segments are then incubated with a large molar excess of linker molecules in the presence of an enzyme that is able to catalyze the ligation of blunt-ended DNA molecules, such as bacteriophage T4 DNA ligase. Thus, the products of the reaction are DNA segments carrying polymeric linker sequences at their ends. These
15 DNA segments are then cleaved with the appropriate restriction enzyme and ligated to an expression vector that has been cleaved with an enzyme that produces termini compatible with those of the DNA segment.

Synthetic linkers containing a variety of restriction endonuclease sites are
20 commercially available from a number of sources including International Biotechnologies Inc, New Haven, CN, USA.

A desirable way to modify the DNA encoding the polypeptide of the invention is to use the polymerase chain reaction as disclosed by Saiki *et*
25 *al* (1988) *Science* 239, 487-491.

In this method the DNA to be enzymatically amplified is flanked by two specific oligonucleotide primers which themselves become incorporated into the amplified DNA. The said specific primers may contain restriction
30 endonuclease recognition sites which can be used for cloning into

expression vectors using methods known in the art.

Variants of the genes also form part of the invention. It is preferred if the variant has at least 70% sequence identity, more preferably at least 85% sequence identity, most preferably at least 95% sequence identity with the genes isolated by the method of the invention. Of course, replacements, deletions and insertions may be tolerated. The degree of similarity between one nucleic acid sequence and another can be determined using the GAP program of the University of Wisconsin Computer Group.

10

Similarly, variants of proteins encoded by the genes are included.

By "variants" we include insertions, deletions and substitutions, either conservative or non-conservative, where such changes do not substantially alter the normal function of the protein.

15

By "conservative substitutions" is intended combinations such as Gly, Ala; Val, Ile, Leu; Asp, Glu; Asn, Gln; Ser, Thr; Lys, Arg; and Phe, Tyr.

Such variants may be made using the well known methods of protein engineering and site-directed mutagenesis.

A ninth aspect of the invention provides a method of identifying a compound which reduces the ability of a microorganism to adapt to a particular environment comprising the steps of selecting a compound which interferes with the function of (1) a gene obtained by the method of the second aspect of the invention or of (2) a polypeptide encoded by such a gene.

25

Pairwise screens for compounds which affect the wild type cell but not a

30

cell overproducing a gene isolated by the methods of the invention form part of this aspect of the invention.

For example, in one embodiment one cell is a wild type cell and a second
5 cell is the *Salmonella* which is made to overexpress the gene isolated by the method of the invention. The viability and/or growth of each cell in the particular environment is determined in the presence of a compound to be tested to identify which compound reduces the viability or growth of the wild type cell but not the cell overexpressing the said gene.

10

It is preferred if the gene is a virulence gene.

For example, in one embodiment the microorganism (such as *S. typhimurium*) is made to over-express the virulence gene identified by the
15 method of the first aspect of the invention. Each of (a) the "over-expressing" microorganism and (b) an equivalent microorganism (which does not over-express the virulence gene) are used to infect cells in culture. Whether a particular test compound will selectively inhibit the virulence gene function is determined by assessing the amount of the test
20 compound which is required to prevent infection of the host cells by (a) the over-expressing microorganism and (b) the equivalent microorganism (at least for some virulence gene products it is envisaged that the test compound will inactivate them, and itself be inactivated, by binding to the virulence gene product). If more of the compound is required to prevent
25 infection by the (a) than (b) then this suggests that the compound is selective. It is preferred if the microorganisms (such as *Salmonella*) are destroyed extracellularly by a mild antibiotic such as gentamicin (which does not penetrate host cells) and that the effect of the test compound in preventing infection of the cell by the microorganism is by lysing the said
30 cell and determining how many microorganisms are present (for example

by plating on agar).

Pairwise screens and other screens for compounds are generally disclosed in Kirsch & Di Domenico (1993) in "The Discovery of Natural Products with a Therapeutic Potential" (Ed, V.P. Gallo), Chapter 6, pages 177-221, 5 Butterworths, V.K. (incorporated herein by reference).

Pairwise screens can be designed in a number of related formats in which the relative sensitivity to a compound is compared using two genetically 10 related strains. If the strains differ at a single locus, then a compound specific for that target can be identified by comparing each strain's sensitivity to the inhibitor. For example, inhibitors specific to the target will be more active against a super-sensitive test strain when compared to an otherwise isogenic sister strain. In an agar diffusion format, this is 15 determined by measuring the size of the zone of inhibition surrounding the disc or well carrying the compound. Because of diffusion, a continuous concentration gradient of compound is set up, and the strain's sensitivity to inhibitors is proportional to the distance from the disc or well to the edge of the zone. General antimicrobials, or antimicrobials with modes 20 of action other than the desired one are generally observed as having similar activities against the two strains.

Another type of molecular genetic screen, involving pairs of strains where a cloned gene product is overexpressed in one strain compared to a control 25 strain. The rationale behind this type of assay is that the strain containing an elevated quantity of the target protein should be more resistant to inhibitors specific to the cloned gene product than an isogenic strain, containing normal amounts of the target protein. In an agar diffusion assay, the zone size surrounding a specific compound is expected to be 30 smaller in the strain overexpressing the target protein compared to an

otherwise isogenic strain.

Additionally or alternatively selection of a compound is achieved in the following steps:

5

1. A mutant microorganism obtained using the method of the first aspect of the invention is used as a control (it has a given phenotype, for example, avirulence).

10 2. A compound to be tested is mixed with the wild-type microorganism.

3. The wild-type microorganism is introduced into the environment (with or without the test compound).

15

4. If the wild-type microorganism is unable to adapt to the environment (following treatment by, or in the presence of, the compound), the compound is one which reduces the ability of the microorganism to adapt to, or survive in, the particular environment.

20

When the environment is an animal body and the microorganism is a pathogenic microorganism, the compound identified by this method can be used in a medicament to prevent or ameliorate infection with the microorganism.

25

A tenth aspect of the invention therefore provides a compound identifiable by the method of the ninth aspect.

It will be appreciated that uses of the compound of the tenth aspect are
30 related to the method by which it can be identified, and in particular in

relation to the host of a pathogenic microorganism. For example, if the compound is identifiable by a method which uses a virulence gene, or polypeptide encoded thereby, from a bacterium which infects a mammal, the compound may be useful in treating infection of a mammal by that
5 bacterium.

Similarly, if the compound is identifiable by a method which uses a virulence gene, or polypeptide encoded thereby, from a fungus which infects a plant, the compound may be useful in treating infection of a plant
10 by that fungus.

An eleventh aspect of the invention provides a molecule which selectively interacts with, and substantially inhibits the function of, a gene of the seventh aspect of the invention or a nucleic acid product thereof.
15

By "nucleic acid product thereof" we include any RNA, especially mRNA, transcribed from the gene.

Preferably a molecule which selectively interacts with, and substantially
20 inhibits the function of, said gene or said nucleic acid product is an antisense nucleic acid or nucleic acid derivative.

More preferably, said molecule is an antisense oligonucleotide.

25 Antisense oligonucleotides are single-stranded nucleic acid, which can specifically bind to a complementary nucleic acid sequence. By binding to the appropriate target sequence, an RNA-RNA, a DNA-DNA, or RNA-DNA duplex is formed. These nucleic acids are often termed "antisense" because they are complementary to the sense or coding strand of the gene.
30 Recently, formation of a triple helix has proven possible where the

oligonucleotide is bound to a DNA duplex. It was found that oligonucleotides could recognise sequences in the major groove of the DNA double helix. A triple helix was formed thereby. This suggests that it is possible to synthesise sequence-specific molecules which specifically
5 bind double-stranded DNA via recognition of major groove hydrogen binding sites.

Clearly, the sequence of the antisense nucleic acid or oligonucleotide can readily be determined by reference to the nucleotide sequence of the gene
10 in question. For example, antisense nucleic acid or oligonucleotides can be designed which are complementary to a part of the sequence shown in Figures 11 or 12, especially to sequences which form a part of a virulence gene.

15 Oligonucleotides are subject to being degraded or inactivated by cellular endogenous nucleases. To counter this problem, it is possible to use modified oligonucleotides, eg having altered internucleotide linkages, in which the naturally occurring phosphodiester linkages have been replaced with another linkage. For example, Agrawal *et al* (1988) *Proc. Natl. Acad. Sci. USA* 85, 7079-7083 showed increased inhibition in tissue culture of HIV-1 using oligonucleotide phosphoramidates and phosphorothioates. Sarin *et al* (1988) *Proc. Natl. Acad. Sci. USA* 85, 7448-7451 demonstrated increased inhibition of HIV-1 using oligonucleotide methylphosphonates. Agrawal *et al* (1989) *Proc. Natl. Acad. Sci. USA* 86, 7790-7794 showed inhibition of HIV-1 replication in both early-infected and chronically infected cell cultures, using nucleotide sequence-specific oligonucleotide phosphorothioates. Leither *et al* (1990) *Proc. Natl. Acad. Sci. USA* 87, 3430-3434 report inhibition in tissue culture of influenza virus replication by oligonucleotide phosphorothioates.

Oligonucleotides having artificial linkages have been shown to be resistant to degradation *in vivo*. For example, Shaw *et al* (1991) in *Nucleic Acids Res.* 19, 747-750, report that otherwise unmodified oligonucleotides become more resistant to nucleases *in vivo* when they are blocked at the
5 3' end by certain capping structures and that uncapped oligonucleotide phosphorothioates are not degraded *in vivo*.

A detailed description of the H-phosphonate approach to synthesizing oligonucleoside phosphorothioates is provided in Agrawal and Tang (1990)
10 *Tetrahedron Letters* 31, 7541-7544, the teachings of which are hereby incorporated herein by reference. Syntheses of oligonucleoside methylphosphonates, phosphorodithioates, phosphoramidates, phosphate esters, bridged phosphoramidates and bridge phosphorothioates are known in the art. See, for example, Agrawal and Goodchild (1987) *Tetrahedron*
15 *Letters* 28, 3539; Nielsen *et al* (1988) *Tetrahedron Letters* 29, 2911; Jager *et al* (1988) *Biochemistry* 27, 7237; Uznanski *et al* (1987) *Tetrahedron Letters* 28, 3401; Bannwarth (1988) *Helv. Chim. Acta.* 71, 1517; Crosstick and Vyle (1989) *Tetrahedron Letters* 30, 4693; Agrawal *et al* (1990) *Proc. Natl. Acad. Sci. USA* 87, 1401-1405, the teachings of which
20 are incorporated herein by reference. Other methods for synthesis or production also are possible. In a preferred embodiment the oligonucleotide is a deoxyribonucleic acid (DNA), although ribonucleic acid (RNA) sequences may also be synthesized and applied.

25 The oligonucleotides useful in the invention preferably are designed to resist degradation by endogenous nucleolytic enzymes. *In vivo* degradation of oligonucleotides produces oligonucleotide breakdown products of reduced length. Such breakdown products are more likely to engage in non-specific hybridization and are less likely to be effective.
30 relative to their full-length counterparts. Thus, it is desirable to use

oligonucleotides that are resistant to degradation in the body and which are able to reach the targeted cells. The present oligonucleotides can be rendered more resistant to degradation *in vivo* by substituting one or more internal artificial internucleotide linkages for the native phosphodiester linkages, for example, by replacing phosphate with sulphur in the linkage. Examples of linkages that may be used include phosphorothioates, methylphosphonates, sulphone, sulphate, ketyl, phosphorodithioates, various phosphoramidates, phosphate esters, bridged phosphorothioates and bridged phosphoramidates. Such examples are illustrative, rather than limiting, since other internucleotide linkages are known in the art. See, for example, Cohen, (1990) *Trends in Biotechnology*. The synthesis of oligonucleotides having one or more of these linkages substituted for the phosphodiester internucleotide linkages is well known in the art, including synthetic pathways for producing oligonucleotides having mixed internucleotide linkages.

Oligonucleotides can be made resistant to extension by endogenous enzymes by "capping" or incorporating similar groups on the 5' or 3' terminal nucleotides. A reagent for capping is commercially available as Amino-Link II™ from Applied BioSystems Inc, Foster City, CA. Methods for capping are described, for example, by Shaw *et al* (1991) *Nucleic Acids Res.* 19, 747-750 and Agrawal *et al* (1991) *Proc. Natl. Acad. Sci. USA* 88(17), 7595-7599, the teachings of which are hereby incorporated herein by reference.

25

A further method of making oligonucleotides resistant to nuclease attack is for them to be "self-stabilized" as described by Tang *et al* (1993) *Nucl. Acids Res.* 21, 2729-2735 incorporated herein by reference. Self-stabilized oligonucleotides have hairpin loop structures at their 3' ends, and show increased resistance to degradation by snake venom

30

phosphodiesterase, DNA polymerase I and fetal bovine serum. The self-stabilized region of the oligonucleotide does not interfere in hybridization with complementary nucleic acids, and pharmacokinetic and stability studies in mice have shown increased *in vivo* persistence of self-stabilized oligonucleotides with respect to their linear counterparts.

In accordance with the invention, the inherent binding specificity of antisense oligonucleotides characteristic of base pairing is enhanced by limiting the availability of the antisense compound to its intended locus *in vivo*, permitting lower dosages to be used and minimizing systemic effects. Thus, oligonucleotides are applied locally to achieve the desired effect. The concentration of the oligonucleotides at the desired locus is much higher than if the oligonucleotides were administered systemically, and the therapeutic effect can be achieved using a significantly lower total amount.

The local high concentration of oligonucleotides enhances penetration of the targeted cells and effectively blocks translation of the target nucleic acid sequences.

The oligonucleotides can be delivered to the locus by any means appropriate for localized administration of a drug. For example, a solution of the oligonucleotides can be injected directly to the site or can be delivered by infusion using an infusion pump. The oligonucleotides also can be incorporated into an implantable device which when placed at the desired site, permits the oligonucleotides to be released into the surrounding locus.

The oligonucleotides are most preferably administered via a hydrogel material. The hydrogel is noninflammatory and biodegradable. Many such materials now are known, including those made from natural and synthetic polymers. In a preferred embodiment, the method exploits a

hydrogel which is liquid below body temperature but gels to form a shape-retaining semisolid hydrogel at or near body temperature. Preferred hydrogel are polymers of ethylene oxide-propylene oxide repeating units. The properties of the polymer are dependent on the molecular weight of the polymer and the relative percentage of polyethylene oxide and polypropylene oxide in the polymer. Preferred hydrogels contain from about 10 to about 80% by weight ethylene oxide and from about 20 to about 90% by weight propylene oxide. A particularly preferred hydrogel contains about 70% polyethylene oxide and 30% polypropylene oxide. Hydrogels which can be used are available, for example, from BASF Corp., Parsippany, NJ, under the tradename Pluronic^R.

In this embodiment, the hydrogel is cooled to a liquid state and the oligonucleotides are admixed into the liquid to a concentration of about 1 mg oligonucleotide per gram of hydrogel. The resulting mixture then is applied onto the surface to be treated, for example by spraying or painting during surgery or using a catheter or endoscopic procedures. As the polymer warms, it solidifies to form a gel, and the oligonucleotides diffuse out of the gel into the surrounding cells over a period of time defined by the exact composition of the gel.

The oligonucleotides can be administered by means of other implants that are commercially available or described in the scientific literature, including liposomes, microcapsules and implantable devices. For example, implants made of biodegradable materials such as polyanhydrides, polyorthoesters, polylactic acid and polyglycolic acid and copolymers thereof, collagen, and protein polymers, or non-biodegradable materials such as ethylenevinyl acetate (EVAc), polyvinyl acetate, ethylene vinyl alcohol, and derivatives thereof can be used to locally deliver the oligonucleotides. The oligonucleotides can be incorporated into the

material as it is polymerized or solidified, using melt or solvent evaporation techniques, or mechanically mixed with the material. In one embodiment, the oligonucleotides are mixed into or applied onto coatings for implantable devices such as dextran coated silica beads, stents, or
5 catheters.

The dose of oligonucleotides is dependent on the size of the oligonucleotides and the purpose for which is it administered. In general, the range is calculated based on the surface area of tissue to be treated.
10 The effective dose of oligonucleotide is somewhat dependent on the length and chemical composition of the oligonucleotide but is generally in the range of about 30 to 3000 μ g per square centimetre of tissue surface area.

The oligonucleotides may be administered to the patient systemically for
15 both therapeutic and prophylactic purposes. The oligonucleotides may be administered by any effective method, for example, parenterally (eg intravenously, subcutaneously, intramuscularly) or by oral, nasal or other means which permit the oligonucleotides to access and circulate in the patient's bloodstream. Oligonucleotides administered systemically
20 preferably are given in addition to locally administered oligonucleotides, but also have utility in the absence of local administration. A dosage in the range of from about 0.1 to about 10 grams per administration to an adult human generally will be effective for this purpose.

25 It will be appreciated that the molecules of this aspect of the invention are useful in treating or preventing any infection caused by the microorganism from which the said gene has been isolated, or a close relative of said microorganism. Thus, the said molecule is an antibiotic.

30 Thus, a twelfth aspect of the invention provides a molecule of the eleventh

aspect of the invention for use in medicine.

A thirteenth aspect of the invention provides a method of treating a host which has, or is susceptible to, an infection with a microorganism, the method comprising administering an effective amount of a molecule according to the eleventh aspect of the invention wherein said gene is present in said microorganisms, or a close relative of said microorganism.

By "effective amount" we mean an amount which substantially prevents or ameliorates the infection. By "host" we include any animal or plant which may be infected by a microorganism.

It will be appreciated that pharmaceutical formulations of the molecule of the eleventh aspect of the invention form part of the invention. Such pharmaceutical formulations comprise the said molecule together with one or more acceptable carriers. The carrier(s) must be "acceptable" in the sense of being compatible with the said molecule of the invention and not deleterious to the recipients thereof. Typically, the carriers will be water or saline which will be sterile and pyrogen free.

20

As mentioned above, and as described in more detail in Example 4 below, I have found that certain virulence genes are clustered in *Salmonella typhimurium* in a region of the chromosome that I have called VGC2. DNA-DNA hybridisation experiments have determined that sequences homologous to at least part of VGC2 are found in many species and strains of *Salmonella* but are not present in the *E. coli* and *Shigella* strains tested (see Example 4). These sequences almost certainly correspond to conserved genes, at least in *Salmonella*, and at least some of which are virulence genes. It is believed that equivalent genes in other *Salmonella* species and, if present, equivalent genes in other enteric or other bacteria

30

will also be virulence genes.

Whether a gene within the VGC2 region is a virulence gene is readily determined. For example, those genes within VGC2 which have been
5 identified by the method of the second aspect of the invention (when applied to *Salmonella typhimurium* and wherein the environment is an animal such as a mouse) are virulence genes. Virulence genes are also identified by making a mutation in the gene (preferably a non-polar mutation) and determining whether the mutant strain is avirulent.
10 Methods of making mutations in a selected gene are well known and are described below.

A fourteenth aspect of the invention provides the VGC2 DNA of *Salmonella typhimurium* or a part thereof, or a variant of said DNA or a
15 variant of a part thereof.

The VGC2 DNA of *Salmonella typhimurium* is depicted diagrammatically in Figure 8 and is readily obtainable from *Salmonella typhimurium* ATCC 14028 (available from the American Type Culture Collection, 12301
20 Parklawn Drive, Rockville, Maryland 20852, USA; also deposited at the NCTC, Public Health Laboratory Service, Colindale, UK under accession no. NCTC 12021) using the information provided in Example 4. For example, probes derived from the sequences shown in Figures 11 and 12 may be used to identify λ clones from a *Salmonella typhimurium* genomic
25 library. Standard genome walking methods can be employed to obtain all of the VGC2 DNA. The restriction map shown in Figure 8 can be used to identify and locate DNA fragments from VGC2.

By "part of the VGC2 DNA of *Salmonella typhimurium*" we mean any
30 DNA sequence which comprises at least 10 nucleotides, preferably at least

20 nucleotides, more preferably at least 50 nucleotides, still more preferably at least 100 nucleotides, and most preferably at least 500 nucleotides of VGC2. A particularly preferred part of the VGC2 DNA is the sequence shown in Figure 11, or a part thereof. Another
5 particularly preferred part of the VGC2 DNA is the sequence shown in Figure 12, or a part thereof.

Advantageously, the part of the VGC2 DNA is a gene, or part thereof.

10 Genes can be identified within the VGC2 region by statistical analysis of the open reading frames using computer programs known in the art. If an open reading frame is greater than about 100 codons it is likely to be a gene (although genes smaller than this are known). Whether an open reading frame corresponds to the polypeptide coding region of a gene can
15 be determined experimentally. For example, a part of the DNA corresponding to the open reading frame may be used as a probe in a northern (RNA) blot to determine whether mRNA is expressed which hybridises to the said DNA; alternatively or additionally a mutation may be introduced into the open reading frame and the effect of the mutation
20 on the phenotype of the microorganism can be determined. If the phenotype is changed then the open reading frame corresponds to a gene. Methods of identifying genes within a DNA sequence are known in the art.

25 By "variant of said DNA or a variant of a part thereof" we include any variant as defined by the term "variant" in the seventh aspect of the invention.

Thus, variants of VGC2 DNA of *Salmonella typhimurium* include
30 equivalent genes, or parts thereof, from other *Salmonella* species, such as

Salmonella typhi and *Salmonella enterica*, as well as equivalent genes, or parts thereof, from other bacteria such as other enteric bacteria.

By "equivalent gene" we include genes which are functionally equivalent and those in which a mutation leads to a similar phenotype (such as avirulence). It will be appreciated that before the present invention VGC2 or the genes contained therein had not been identified and certainly not implicated in virulence determination.

Thus, further aspects of the invention provide a mutant bacterium wherein if the bacterium normally contains a gene that is the same as or equivalent to a gene in VGC2, said gene is mutated or absent in said mutant bacterium; methods of making a mutant bacterium wherein if the bacterium normally contains a gene that is the same as or equivalent to a gene in VGC2, said gene is mutated or absent in said mutant bacterium.

The following is a preferred method to inactivate a VGC2 gene. One first subclones the gene on a DNA fragment from a *Salmonella* λ DNA library or other DNA library using a fragment of VGC2 as a probe in hybridisation experiments, and map the gene with respect to restriction enzyme sites and characterise the gene by DNA sequencing in *Escherichia coli*. Using restriction enzymes, one then introduces into the coding region of the gene a segment of DNA encoding resistance to an antibiotic (for example, kanamycin), possibly after deleting a portion of the coding region of the cloned gene by restriction enzymes. Methods and DNA constructs containing an antibiotic resistance marker are available to ensure that the inactivation of the gene of interest is preferably non-polar, that is to say, does not affect the expression of genes downstream from the gene of interest. The mutant version of the gene is then transferred from *E. coli* to *Salmonella typhimurium* using phage P22 transduction and transductants checked by Southern hybridisation for homologous

recombination of the mutant gene into the chromosome.

This approach is commonly used in *Salmonella* (and can be used in *S. typhi*), and further details can be found in many papers, including Galan
5 *et al* (1992) 174, 4338-4349.

Still further aspects provide a use of said mutant mutant bacterium in a vaccine; pharmaceutical compositions comprising said bacterium and a pharmaceutically acceptable carrier; a polypeptide encoded by VGC2
10 DNA of *Salmonella typhimurium* or a part thereof, or a variant of a part thereof; a method of identifying a compound which reduces the ability of a bacterium to infect or cause disease in a host; a compound identifiable by said method; a molecule which selectively interacts with, and substantially inhibits the function of, a gene in VGC2 or a nucleic product
15 thereof; and medical uses and pharmaceutical compositions thereof.

The VGC2 DNA contains genes which have been identified by the methods of the first and second aspects of the invention as well as genes which have been identified by their location (although identifiable by the
20 methods of the first and second aspects of the invention). These further aspects of the invention relate closely to the fourth, fifth, sixth, seventh, eighth, ninth, tenth, eleventh, twelfth and thirteenth aspects of the invention and, accordingly, the information given in relation to those aspects, and preferences expressed in relation to those aspects, applies to
25 these further aspects.

It is preferred if the gene is from VGC2 or is an equivalent gene from another species of *Salmonella* such as *S. typhi*. It is preferred if the mutant bacterium is a *S. typhimurium* mutant or a mutant of another
30 species of *Salmonella* such as *S. typhi*.

It is believed that at least some of the genes in VGC2 confer the ability for the bacterium, such as *S. typhimurium*, to enter cells.

The invention will now be described with reference to the following
5 Examples and Figures wherein:

Figure 1 illustrates diagrammatically one particularly preferred method of the invention.

10 Figure 2 shows a Southern hybridisation analysis of DNA from 12 *S. typhimurium* exconjugants following digestion with *EcoRV*. The filter was probed with the kanamycin resistance gene of the mini-Tn5 transposon.

Figure 3 shows a colony blot hybridisation analysis of DNA from 48 *S.*
15 *typhimurium* exconjugants from a half of a microtitre dish (A1-H6). The filter was hybridised with a probe comprising labelled amplified tags from DNA isolated from a pool of the first 24 colonies (A1-D6).

Figure 4 shows a DNA colony blot hybridisation analysis of 95 *S.*
20 *typhimurium* exconjugants of a microtitre dish (A1-H11), which were injected into a mouse. Replicate filters were hybridised with labelled amplified tags from the pool (inoculum pattern), or with labelled amplified tags from DNA isolated from over 10,000 pooled colonies that were recovered from the spleen of the infected animal (spleen pattern).
25 Colonies B6, A11 and C8 gave rise to weak hybridisation signals on both sets of filters. Hybridisation signals from colonies A3, C5, G3 (*aroA*), and F10 are present on the inoculum pattern but not on the spleen pattern.

Figure 5 shows the sequence of a *Salmonella* gene isolated using the
30 method of the invention and a comparison to the *Escherichia coli* *clp*

protease genome.

Figure 6 shows partial sequences of further *Salmonella* gene isolated using the method of the invention (SEQ ID Nos. 8 to 36).

5

Figure 7 shows the mapping of VGC2 on the *S. typhimurium* chromosome. (A) DNA probes from three regions of VGC2 were used in Southern hybridisation analysis of lysates from a set of *S. typhimurium* strains harbouring locked in Mud-P22 prophages. Lysates which
10 hybridised to a 7.5 kb *Pst*I fragment (probe A in Figure 8) are shown. The other two probes used hybridised to the same lysates. (B) The insertion points and packaging directions of the phage are shown along with the map position in minutes (edition VIII, ref 22 in Example 4). The phage designations correspond to the following strains: 18P, TT15242;
15 18Q, 15241; 19P, TT15244; 19Q, TT15243; 20P, TT15246 and 20Q, TT15245 (Ref in Example 4). The locations of mapped genes are shown by horizontal bars and the approximate locations of other genes are indicated.

20 Figure 8 shows a physical and genetic map of VGC2. (A) The positions of 16 transposon insertions are shown above the line. The extent of VGC2 is indicated by the thicker line. The position and direction of transcription of ORFs described in the text of Example 4 are shown by arrows below the line, together with the names of similar genes, with the
25 exception of ORFs 12 and 13 whose products are similar to the sensor and regulatory components respectively, of a variety of two component regulatory systems. (B) The location of overlapping clones and an *Eco*RI/*Xba*I restriction fragment from Mud-P22 prophage strain TT15244 are shown as filled bars. Only the portions of the λ clones which have
30 been mapped are shown and the clones may extend beyond these limits.

(C) The positions of restriction sites are marked: B, *Bam*HI; E, *Eco*RI; V, *Eco*RV; H, *Hind*III; P, *Pst*I and X, *Xba*I. The positions of the 7.5 kb *Pst*I fragment (probe A) used as a probe in Figure 7 and that of the 2.2 kb *Pst*I/*Hind*III fragment (probe B) used as a probe in Figure 10 are shown below the restriction map. The positions of Sequence 1 (described in Figure 11) and Sequence 2 (described in Figure 12) are shown by the thin arrows (labelled Sequence 1 and Sequence 2).

Figure 9 describes mapping the boundaries of VGC2. (A) The positions of mapped genes at minutes 37 to 38 on the *E. coli* K12 chromosome are aligned with the corresponding region of the *S. typhimurium* LT2 chromosome (minutes 30 to 31). An expanded map of the VGC2 region is shown with 11 *S. typhimurium* (*S. t.*) DNA fragments used as probes (thick bars) and the restriction sites used to generate them: B, *Bam*HI; C, *Cla*I; H, *Hind*II; K, *Kpn*I; P, *Pst*I; N, *Nsi*I and S, *Sal*I. Probes that hybridised to *E. coli* K12 (*E. c.*) genomic DNA are indicated by +; those which failed to hybridise are indicated by -.

Figure 10 shows that VGC2 is conserved among and specific to the *Salmonellae*. Genomic DNA from *Salmonella* serovars and other pathogenic bacteria was restricted with *Pst*I (A), *Hind*III or *Eco*RV (B) and subjected to Southern hybridisation analysis, using a 2.2 kb *Pst*I/*Hind*III fragment from λ clone 7 as a probe (probe B Figure 2). The filters were hybridised and washed under stringent (A) or non-stringent (B) conditions.

Figure 11 shows the DNA sequence of "Sequence 1" of VGC2 from the centre to the left-hand end (see the arrow labelled Sequence 1 in Figure 2). The DNA is translated in all six reading frames and the start and stop positions of putative genes, and the transposon insertion positions for

various mutants identified by STM are indicated (SEQ ID No 37).

As is conventional a * indicates a stop codon and standard nucleotide ambiguity codes are used where necessary.

5

Figure 12 shows the DNA sequence of "Sequence 2" of VGC2 (cluster C) (see the arrow labelled Sequence 2 in Figure 2). The DNA is translated in all six reading frames and the start and stop positions of putative genes, and the transposon insertion positions for various mutants identified by STM are indicated (SEQ ID No 38).

10

As is conventional a * indicates a stop codon and standard nucleotide ambiguity codes are used where necessary.

15 Figures 7 to 12 are most relevant to Example 4.

Example 1: Identification of virulence genes in *Salmonella typhimurium*

20 **Materials and Methods**

Bacterial Strains and Plasmids

Salmonella typhimurium strain 12023 (equivalent to American Type Culture Collection (ATCC) strain 14028) was obtained from the National Collection of Type Cultures (NCTC), Public Health Laboratory Service, Colindale, London, UK. A spontaneous nalidixic acid resistant mutant of this strain (12023 Nal^r) was selected in our laboratory. Another derivative of strain 12023, CL1509 (*aroA::Tn10*) was a gift from Fred Heffron.

25

30 *Escherichia coli* strains CC118 λ pir (Δ [*ara-leu*], *araD*, Δ *lacX74*, *galE*,

galK, *phoA20*, *thi-1*, *rpsE*, *rpoB*, *argE*(Am), *recA1*, λ pir phage lysogen and S17-1 λ pir (Tp^r, Sm^r, *recA*, *thi*, *pro*, *hsdR*⁻M⁺, RP4:2-Tc:Mu:KmTn7, λ pir) were gifts from Kenneth Timmis. *E. coli* DH5 α was used for propagating pUC18 (Gibco-BRL) and Bluescript (Stratagene) plasmids
5 containing *S. typhimurium* DNA. Plasmid pUTmini-Tn5Km2 (de Lorenzo *et al*, 1990) was a gift from Kenneth Timmis.

Construction of semi-random sequence tags and ligations

10 The oligonucleotide pool RT1(5'-CTAGGTACCTACAACCTCAAGCTT-[NK]₂₀-AAGCTTGGTTAGAATGGGTACCATG-3') (SEQ ID No 1), and primers P2 (5'-TACCTACAACCTCAAGCT-3') (SEQ ID No 2), P3 (5'-CATGGTACCCATTCTAAC-3') (SEQ ID No 3), P4 (5'-TACCCATTCTAACCAAGC-3') (SEQ ID No 4) and P5 (5'-
15 CTAGGTACCTACAACCTC-3') (SEQ ID No 5) were synthesized on a oligonucleotide synthesizer (Applied Biosystems, model 380B). Double stranded DNA tags were prepared from RT1 in a 100 μ l volume PCR containing 1.5 mM MgCl₂, 50 mM KCl, and 10 mM Tris-Cl (pH 8.0) with 200 pg of RT1 as target; 250 μ M each dATP, dCTP, dGTP, dTTP; 100 pM of
20 primers P3 and P5; and 2.5 U of Amplitaq (Perkin-Elmer Cetus). Thermal cycling conditions were 30 cycles of 95°C for 30 s, 50°C for 45 s, and 72°C for 10 s. The PCR product was gel purified (Sambrook *et al*, 1989), passed through an elutipD column (available from Schleicher and Schull) and digested with *KpnI* prior to ligation into pUC18 or pUTmini-Tn5Km2. For ligations,
25 plasmids were digested with *KpnI* and dephosphorylated with calf intestinal alkaline phosphatase (Gibco-BRL). Linearized plasmid molecules were gel-purified (Sambrook *et al*, 1989) prior to ligation to remove any residual uncut plasmid DNA from the digestion. Ligation reactions contained approximately 50 ng each of plasmid and double stranded tag DNA in a 25 μ l volume with 1
30 unit T4 DNA ligase (Gibco-BRL) in a buffer supplied with the enzyme.

Ligations were carried out for 2 h at 24°C. To determine the proportion of bacterial colonies arising from either self ligation of the plasmid DNA or uncut plasmid DNA, a control reaction was carried out in which the double stranded tag DNA was omitted from the ligation reaction. This yielded no ampicillin resistant bacterial colonies following transformation of *E. coli* CC118 (Sambrook *et al*, 1989), compared with 185 colonies arising from a ligation reaction containing the double stranded tag DNA.

Bacterial Transformation and Matings

10

The products of several ligations between pUT mini-Tn5Km2 and the double stranded tag DNA were used to transform *E. coli* CC118 (Sambrook *et al*, 1989). A total of approximately 10,300 transformants were pooled and plasmid DNA extracted from the pool was used to transform *E. coli* S-17 λ pir (de Lorenzo & Timmis, 1994). For mating experiments, a pool of approximately 40,000 ampicillin resistant *E. coli* S-17 λ pir transformants, and *S. typhimurium* 12023 Nal^r were cultured separately to an optical density (OD)₅₈₀ of 1.0. Aliquots of each culture (0.4 ml) were mixed in 5 ml 10 mM MgSO₄, and filtered through a Millipore membrane (0.45 μ m diameter). The filters were placed on the surface of agar containing M9 salts (de Lorenzo & Timmis, 1994) and incubated at 37°C for 16 h. The bacteria were recovered by shaking the filters in liquid LB medium for 40 min at 37°C and exconjugants were selected by plating the suspension onto LB medium containing 100 μ g ml⁻¹ nalidixic acid (to select against the donor strain) and 50 μ g ml⁻¹ kanamycin (to select for the recipient strain). Each exconjugant was checked by transferring nalidixic acid resistant (nal^r), kanamycin resistant (kan^r) colonies to MacConkey Lactose indicator medium (to distinguish between *E. coli* and *S. typhimurium*), and to LB medium containing ampicillin. Approximately 90% of the nal^r, kan^r colonies were sensitive to ampicillin, indicating that these resulted from authentic

30

transposition events (de Lorenzo & Timmis, 1994). Individual ampicillin-sensitive exconjugants were stored in 96 well microtitre dishes containing LB medium. For long term storage at -80°C, either 7% DMSO or 15% glycerol was included in the medium.

5

Phenotypic characterisation of mutants

Mutants were replica plated from microtitre dishes onto solid medium containing M9 salts and 0.4% glucose (Sambrook *et al*, 1989) to identify
10 auxotrophs. Mutants with rough colony morphology were detected by low magnification microscopy of colonies on agar plates.

Colony Blots, DNA extractions, PCRs, DNA labelings and hybridisations

15 For colony blot hybridizations, a 48-well metal replicator (Sigma) was used to transfer exconjugants from microtitre dishes to Hybond N nylon filters (Amersham, UK) that had been placed on the surface of LB agar containing 50 µg ml⁻¹ kanamycin. After overnight incubation at 37°C, the filters supporting the bacterial colonies were removed and dried at room
20 temperature for 10 min. The bacteria were lysed with 0.4 N NaOH and the filters washed with 0.5 N Tris-Cl pH 7.0 according to the filter manufacturer's instructions. The bacterial DNA was fixed to the filters by exposure to UV light from a Stratalinker (Stratagene). Hybridisations to ³²P-labelled probes were carried out under stringent conditions as previously
25 described (Holden *et al*, 1989). For DNA extractions, *S. typhimurium* transposon mutant strains were grown in liquid LB medium in microtitre dishes or resuspended in LB medium following growth on solid media. Total DNA was prepared by the hexadecyltrimethylammoniumbromide (CTAB) method according to Ausubel *et al* (1987). Briefly, cells from 150
30 to 1000 µl volumes were precipitated by centrifugation and resuspended in

- 576 μ l TE. To this was added 15 μ l of 20% SDS and 3 μ l of 20 mg ml⁻¹ proteinase K. After incubating at 37°C for 1 hour, 166 μ l of 3 M NaCl was added and mixed thoroughly, followed by 80 μ l of 10% (w/v) CTAB and 0.7 M NaCl. After thorough mixing, the solution was incubated at 5 65°C for 10 min. Following extraction with phenol and phenol-chloroform, the DNA was precipitated by addition of isopropanol, washed with 70% ethanol and resuspended in TE at a concentration of approximately 1 μ g μ l⁻¹.
- 10 The DNA samples were subjected to two rounds of PCR to generate labelled probes. The first PCR was performed in 100 μ l reactions containing 20 mM Tris-Cl pH 8.3; 50 mM KCl; 2 mM MgCl₂; 0.01% Tween 80; 200 μ M each dATP, dCTP, dGTP, dTTP; 2.5 units of Amplitaq polymerase (Perkin-Elmer Cetus); 770 ng each primer P2 and P4; and 5 μ g target DNA. After an initial denaturation of 4 min at 95°C, thermal cycling consisted of 20 cycles of 45 s at 50°C, 10 s at 72°C, and 30 s at 95°C. PCR products were extracted with chloroform/isoamyl alcohol (24/1) and precipitated with ethanol. DNA was resuspended in 10 μ l TE and the PCR products were purified by electrophoresis through a 1.6% Seaplaque (FMC 20 Bioproducts) gel in TAE buffer. Gel slices containing fragments of about 80 bp were excised and used for the second PCR. This reaction was carried out in a 20 μ l total volume, and contained 20 mM Tris-Cl pH 8.3; 50 mM KCl; 2 mM MgCl₂; 0.01% Tween 80; 50 μ M each dATP, dTTP, dGTP; 10 μ l ³²P-dCTP (3000 Ci/mmol, Amersham); 150 ng each primer P2 and P4; approximately 10 ng of target DNA (1-2 μ l of 1.6% Seaplaque 25 agarose containing the first round PCR product); 0.5 units of Amplitaq polymerase. The reaction was overlayed with 20 μ l mineral oil and thermal cycling was performed as described above. Incorporation of the radioactive label was quantitated by absorbance to Whatman DE81 paper (Sambrook *et al*, 1989). 30

Infection Studies

Individual *Salmonella* exconjugants containing tagged transposons were grown in 2% tryptone, 1% yeast extract, 0.92% v/v glycerol, 0.5% Na₂PO₄, 1% KNO₃ (TYGPN medium) (Ausubel *et al*, 1987) in microtitre plates overnight at 37°C. A metal replicator was used to transfer a small volume of the overnight cultures to a fresh microtitre plate and the cultures were incubated at 37°C until the OD₅₈₀ (measured using a Titertek Multiscan microtitre plate reader) was approximately 0.2 in each well. Cultures from individual wells were then pooled and the OD₅₈₀ determined using a spectrophotometer. The culture was diluted in sterile saline to approximately 5x10⁵ cfu ml⁻¹. Further dilutions were plated out onto TYGPN containing nalidixic acid (100 mg ml⁻¹) and kanamycin (50 mg ml⁻¹) to confirm the cfu present in the inoculum.

Groups of three female BALB/c mice (20-25g) were injected intraperitoneally with 0.2 ml of bacterial suspension containing approximately 1x10⁵ cfu ml⁻¹. Mice were sacrificed three days post-inoculation and their spleens were removed to recover bacteria. Half of each spleen was homogenized in 1 ml of sterile saline in a microfuge tube. Cellular debris was allowed to settle and 1 ml of saline containing cells still in suspension was removed to a fresh tube and centrifuged for two minutes in a microfuge. The supernatant was aspirated and the pellet resuspended in 1 ml of sterile distilled water. A dilution series was made in sterile distilled water and 100 ml of each dilution was plated onto TYGPN agar containing nalidixic acid (100 ug ml⁻¹) and kanamycin (50 ug ml⁻¹). Bacteria were recovered from plates containing between 1000 and 4000 colonies, and a total of over 10,000 colonies recovered from each spleen were pooled and used to prepare DNA for PCR generation of probes to screen colony blots.

Virulence gene cloning and DNA sequencing

Total DNA was isolated from *S. typhimurium* exconjugants and digested separately with *Sst*I, *Sal*I, *Pst*I and *Sph*I. Digests were fractionated through agarose gels, transferred to Hybond N⁺ membranes (Amersham) and subjected to Southern hybridisation analysis using the kanamycin resistance gene of pUT mini-Tn5Km2 as a probe. The probe was labelled with digoxigenin (Boehringer-Mannheim) and chemiluminescence detection was carried out according to the manufacturer's instructions. The hybridisation and washing conditions were as described above. Restriction enzymes which gave rise to hybridising fragments in the 3-5 kb range were used to digest DNA for a preparative agarose gel, and DNA fragments corresponding to the sizes of the hybridisation signals were excised from this, purified and ligated into pUC18. Ligation reactions were used to transform *E. coli* DH5a to kanamycin resistance. Plasmids from kanamycin-resistant transformants were purified by passage through an elutipD column and checked by restriction enzyme digestion. Plasmid inserts were partially sequenced by the di-deoxy method (Sanger *et al*, 1977) using the -40 primer and reverse sequencing primer (United States Biochemical Corporation) and the primers P6 (5'-CCTAGGCGGCCAGATCTGAT-3') (SEQ ID No 6) and P7 (5'-GCACTTGTGTATAAGAGTCAG-3') (SEQ ID No 7) which anneal to the I and O termini of Tn5, respectively. Nucleotide sequences and deduced amino acid sequences were assembled using the Macvector 3.5 software package run on a Macintosh SE/30 computer. Sequences were compared with the EMBL and Genbank DNA databases using the UNIX/SUN computer system at the Human Genome Mapping Project Resource Centre, Harrow, UK.

Results

Tag Design

5 The structure of the DNA tags is shown in Figure 1a. Each tag consists of a variable central region flanked by "arms", of invariant sequence. The central region sequence ([NK]₂₀) was designed to prevent the occurrence of sites for the commonly used 6 bp-recognition restriction enzymes, but is sufficiently variable to ensure that statistically, the same sequence should
10 only occur once in 2×10^{11} molecules (DNA sequencing of 12 randomly selected tags showed that none shared more than 50% identity over the variable region). (N means any base (A, G, C or T) and K means G or T.) The arms contain *KpnI* sites close to the ends to facilitate the initial cloning step, and the *HindIII* sites bordering the variable region were used to release
15 radiolabelled variable regions from the arms prior to hybridisation analysis. The arms were also designed such that primers P2 and P4 each contain only one guanine residue. Therefore during a PCR using these primers, only one cytosine will be incorporated into each newly synthesised arm, compared to an average of ten in the unique sequence. When radiolabelled dCTP is
20 included in the PCR, an average of ten-fold more label will be present in the unique sequence compared with each arm. This is intended to minimise background hybridisation signals from the arms, after they have been released from the unique sequences by digestion with *HindIII*. Double stranded tags were ligated into the *KpnI* site of the mini-Tn5 transposon
25 Km2, carried on plasmid pUT (de Lorenzo & Timmis, 1994). Replication of this plasmid is dependent on the R6K-specified π product of the *pir* gene. It carries the *oriT* sequence of the RP4 plasmid, permitting transfer to a variety of bacterial species (Miller & Mekalanos, 1988), and the *tnp*⁺ gene needed for transposition of the mini-Tn5 element. The tagged mini-Tn5
30 transposons were transferred to *S. typhimurium* by conjugation, and 288

exconjugants resulting from transposition events were stored in the wells of microtitre dishes. Total DNA isolated from 12 of these was digested with *EcoRV*, and subjected to Southern hybridisation analysis using the kanamycin resistance gene of the mini-Tn5 transposon as a probe. In each case, the exconjugant had arisen as a result a single integration of the transposon into a different site of the bacterial genome (Figure 2).

Specificity and sensitivity studies

We next determined the efficiency and uniformity of amplification of the DNA tags in PCRs involving pools of exconjugant DNAs as targets for the reactions. In an attempt to minimise unequal amplification of tags in the PCR, we determined the maximum quantity of DNA target that could be used in a 100 μ l reaction, and the minimum number of PCR cycles, that resulted in products which could be visualised by ethidium bromide staining of an agarose gel (5 μ g DNA and 20 cycles, respectively).

S. typhimurium exconjugants which had reached stationary growth phase in microtitre dishes were combined, and used to extract DNA. This was subjected to a PCR using primers P2 and P4. PCR products of 80 bp were gel-purified and used as targets for a second PCR, using the same primers but with 32 P-labelled CTP. This resulted in over 60% of the radiolabelled dCTP being incorporated into the PCR products. The radiolabelled products were digested with *HindIII* and used to probe colony blotted DNA from their corresponding microtitre dishes. Of the 1510 mutants tested in this way, 358 failed to yield a clear signal on an autoradiogram following an overnight exposure of the colony blot. There are three potential explanations for this. Firstly, it is possible that a proportion of the transposons did not carry tags. However, by comparing the transformation frequencies resulting from ligation reactions involving the transposon in the

presence or absence of tags, it seems unlikely that untagged transposons could account for more than approximately 0.5% of the total (see Materials and Methods). More probable causes are that the variable sequence was truncated in some of the tags, and/or that some of the sequences formed
5 secondary structures, both of which might have prevented amplification. Mutants which failed to give clear signals were not included in further studies. The specificity of the efficiently amplifiable tags was demonstrated by generating a probe from 24 colonies of a microtitre dish, and using it to probe a colony blot of 48 colonies, which included the 24 used to generate
10 the probe. The lack of any hybridisation signal from the 24 colonies not used to generate the probe (Figure 3) shows that the hybridisation conditions employed were sufficiently stringent to prevent cross-hybridisation among labelled tags, and suggests that each exconjugant is not reiterated within a microtitre dish.

15

There are further considerations in determining the maximum pool size that can be used as an inoculum in animal experiments. As the quantity of labelled tag for each transposon is inversely proportional to the complexity of the tag pool, there is a limit to the pool size above which hybridisation
20 signals become too weak to be detected after overnight exposure of an autoradiogram. More importantly, as the complexity of the pool increases, so must the likelihood of failure of a virulent representative of the pool to be present in sufficient numbers, in the spleen of an infected animal, to produce enough labelled probe. We have not determined the upper limit for
25 pool size in the murine model of salmonellosis that we have employed, but it must be in excess of 96.

Virulence tests of the transposon mutants

30 A total of 1152 uniquely tagged insertion mutants (from two microtitre

dishes) were tested for virulence in BALB/c mice in twelve pools, each representing a 96-well microtitre dish. Animals received an intraperitoneal injection of approximately 10^3 cells of each of 96 transposon mutants of a microtitre dish (10^5 organisms in total). Three days after injection mice were sacrificed, and bacteria were recovered by plating spleen homogenates onto laboratory medium. Approximately 10,000 colonies recovered from each mouse were pooled and DNA was extracted. The tags present in this DNA sample were amplified and labelled by the PCR, and colony blots probed and compared with the hybridisation pattern obtained using tags amplified from the inoculum (Figure 3). As a control, an *aroA* mutant of *S. typhimurium* was tagged and employed as one of the 96 mutants in the inoculum. This strain would not be expected to be recovered in the spleen because its virulence is severely attenuated (Buchmeier *et al*, 1993). Forty-one mutants were identified whose DNA hybridized to labelled tags from the inoculum but not from labelled tags from bacteria recovered from the spleen. The experiment was repeated and the same forty-one mutants were again identified. Two of these were the *aroA* mutant (one per pool), as expected. Another was an auxotrophic mutant (it failed to grow on minimal medium). All of the mutants had normal colony morphology.

20

Example 2: Cloning and partial characterisation of sequences flanking the transposon

DNA was extracted from one of the mutants described in Example 1 (Pool 1, F10), digested with *Sst*I, and subcloned on the basis of kanamycin resistance. The sequence of 450 bp flanking one end of the transposon was determined using primer P7. This sequence shows 80% identity to the *E. coli* *clp* (*lon*) gene, which encodes a heat-regulated protease (Figure 5). To our knowledge, this gene has not previously been implicated as a virulence determinant.

30

Partial sequences of thirteen further *Salmonella typhimurium* virulence genes are shown in Figure 6 (sequences A2 to A9 and B1 to B5). Deduced amino acid sequences of P2D6, S4C3, P3F4, P7G2 and P9B7 bear similarities to a family of secretion-associated proteins that have been conserved throughout bacterial pathogens of animals and plants, and which are known in *Salmonella* as the *inv* family. In *S. typhimurium* the *inv* genes are required for bacterial invasion into intestinal tissue. The virulence of *inv* mutants is attenuated when they are inoculated by the oral route, but not when they are administered intraperitoneally. The discovery of *inv*-related genes that are required for virulence following intraperitoneal inoculation suggests a new secretion apparatus which might be required for invasion of non-phagocytic cells of the spleen and other organs. The products of these new genes might represent better drug targets than the *inv* proteins in the treatment of established infections.

15

Further characterisation of the genes identified in this example is described in Example 4.

Example 3: LD₅₀ determinations and mouse vaccination study

20

Mutations identified by the method of the invention attenuate virulence.

Five of the mutations in genes not previously implicated in virulence were transferred by P22-mediated transduction to the nalidixic acid-sensitive parent strain of *S. typhimurium* 12028. Transductants were checked by restriction mapping then injected by the intraperitoneal route into groups of BALB/c mice to determine their 50% lethal dose (LD₅₀). The LD₅₀ values for mutants S4C3, P7G2, P3F4 and P9B7 were all several orders of magnitude higher than that of the wild-type strain. No difference in the LD₅₀ was detected for mutant P1F10; however, there was a statistically

30

significant decrease in the proportion of P1F10 cells recovered from the spleens of mice injected with an inoculum consisting of an equal proportion of this strain and the wild-type strain. This implies that this mutation does attenuate virulence, but to a degree that is not detectable by LD₅₀.

5

Mutants P3F4 and P9B7 were also administered by the oral route at an inoculum level of 10⁷ cells/mouse. None of the mice became ill, indicating that the oral LD₅₀ levels of these mutants are at least an order of magnitude higher than that of the wild-type strain.

10

In the mouse vaccination study groups of five female BALB/c mice of 20-25 g in mass were initially inoculated orally (p.o.) or intraperitoneally (i.p.) with serial ten fold dilutions of *Salmonella typhimurium* mutant strains P3F4 and P9B7. After four weeks the mice were then inoculated with 500 c.f.u.

15 of the parental wild type strain. Deaths were then recorded over four weeks.

A group of two mice of the same age and batch as the mice inoculated with the mutant strains were also inoculated i.p. with 500 c.f.u. of the wild type strain as a positive control. Both non-immunised mice died as expected within four weeks.

20

Results are tabulated below:

25 1) p.o. initial inoculation with mutant strain P3F4

initial inoculum in c.f.u.	no. mice surviving first challenge	no. mice surviving wild type challenge
5 x 10 ⁹	5	2 (40%)
5 x 10 ⁸	5	2 (40%)

5×10^7	5	0 (0%)
-----------------	---	--------

2) i.p. initial inoculum with mutant strain P3F4

5	initial inoculum in c.f.u.	no. mice surviving first challenge	no. mice surviving wild type challenge
	5×10^6	3	3 (100%)
	5×10^5	5	4 (80%)
	5×10^4	6	5 (83%)
10	5×10^3	5	4 (80%)

3) p.o. initial inoculum with mutant strain P9B7

15	initial inoculum in c.f.u.	no. mice surviving first challenge	no. mice surviving wild type challenge
	5×10^9	5	0 (0%)

4) i.p. initial inoculum with mutant P9B7

20	initial inoculum in c.f.u.	no. mice surviving first challenge	no. mice surviving wild type challenge
	5×10^6	4	2 (50%)

From these experiments I conclude that mutant P3P4 appears to give some protection against subsequent wild type challenge. This protection appears greater in mice that were immunised i.p.

Example 4: Identification of a virulence locus encoding a second type III secretion system in *Salmonella typhimurium*

30

Abbreviations used in this Example are VGC1, virulence gene cluster 1; VGC2, virulence gene cluster 2.

Background to the experiments described

Salmonella typhimurium is a principal agent of gastroenteritis in humans and produces a systemic illness in mice which serves as a model for human typhoid fever (1). Following oral inoculation of mice with *S. typhimurium*, the bacteria pass from the lumen of the small intestine through the intestinal mucosa, via enterocytes or M cells of the Peyer's patch follicles (2). The bacteria then invade macrophages and neutrophils, enter the reticuloendothelial system and disseminate to other organs, including the spleen and liver, where further reproduction results in an overwhelming and fatal bacteremia (3). To invade host cells, to survive and replicate in a variety of physiologically stressful intracellular and extracellular environments and to circumvent the specific antibacterial activities of the immune system, *S. typhimurium* employs a sophisticated repertoire of virulence factors (4).

To gain a more comprehensive understanding of virulence mechanisms of *S. typhimurium* and other pathogens the transposon mutagenesis system described in Example 1, which is conveniently called 'signature-tagged mutagenesis' (STM), which combines the strength of mutational analysis with the ability to follow simultaneously the fate of a large number of different mutants within a single animal (5 and Example 1; Reference 5 was published after the priority date for this invention). Using this approach we identified 43 mutants with attenuated virulence from a total of 1152 mutants that were screened. The nucleotide sequences of DNA flanking the insertion points of transposons in 5 of these mutants showed that they were related to genes encoding type III secretion systems of a variety of bacterial pathogens (6, 7). The products of the *inv/spa* gene cluster of *S. typhimurium* (8, 9) are proteins that form a type III secretion system required for the assembly of surface appendages mediating entry into

epithelial cells (10). Hence the virulence of strains carrying mutations in the *inv/spa* cluster is attenuated only if the inoculum is administered orally and not when given intraperitoneally (8). In contrast the 5 mutants identified by STM are avirulent following intraperitoneal inoculation (5).

5

In this example we show that the transposon insertion points of these 5 mutants and an additional 11 mutants identified by STM all map to the same region of the *S. typhimurium* chromosome. Further analysis of this region reveals additional genes whose deduced products have sequence similarity to other components of type III secretion systems. This chromosomal region which we refer to as virulence gene cluster 2 (VGC2) is not present in a number of other enteric bacteria, and represents an important locus for *S. typhimurium* virulence.

15 Materials and Methods

Bacterial Strains, Transduction and Growth Media. *Salmonella enterica* serotypes 5791 (*abderdeen*), 423180 (*gallinarum*), 7101 (*cubana*) and 12416 (*typhimurium* LT2) were obtained from the National Collections of Type Cultures, Public Health Laboratory Service, UK. *Salmonella typhi* BRD123 genomic DNA was a gift from G. Dougan, enteropathogenic *Escherichia coli* (EPEC), enterohemorrhagic *E. coli* (EHEC), *Vibrio cholera* biotype *El Tor*, *Shigella flexneri* serotype 2 and *Staphylococcus aureus* were clinical isolates obtained from the Department of Infectious Diseases and Bacteriology, Royal Postgraduate Medical School, UK. Genomic DNA from *Yersinia pestis* was a gift from J. Heesemann. However, genomic DNA can be isolated using standard methods. The bacterial strains and the methods used to generate signature-tagged mini-Tn5 transposon mutants of *S. typhimurium* NCTC strain 12023 have been described previously (5, 11).

30 Routine propagation of plasmids was in *E. coli* DH5 α . Bacteria were

grown in LB broth (12) supplemented with the appropriate antibiotics. Before virulence levels of individual mutant strains were assessed, the mutations were first transferred by phage P22 mediated transduction (12) to the nalidixic acid sensitive parental strain of *S. typhimurium* 12023.

- 5 Transductants were analysed by restriction digestion and Southern hybridisation before use as inoculum.

Lambda Library Screening. Lambda (λ) clones with overlapping insert DNAs covering VGC2 were obtained by standard methods (13) from a
10 λ 1059 library (14) containing inserts from a partial *Sau*3A digest of *S. typhimurium* LT2 genomic DNA. The library was obtained via K. Sanderson, from the Salmonella Genetic Stock Centre (SGSC), Calgary, Canada.

- 15 **Mud-P22 Lysogens.** Radiolabelled DNA probes were hybridised to Hybond N (Amersham) filters bearing DNA prepared from lysates of a set of *S. typhimurium* strains harbouring Mud-P22 prophages at known positions in the *S. typhimurium* genome. Preparation of mitomycin-induced Mud-P22 lysates was as described (12, 15). The set of Mud-P22 prophages
20 was originally assembled by Benson and Goldman (16) and was obtained from the SGSC.

Gel Electrophoresis and Southern Hybridisation. Gel electrophoresis was performed in 1% or 0.6% agarose gels run in 0.5 x TBE. Gel fractionated
25 DNA was transferred to Hybond N or N+ membranes (Amersham) and stringent hybridisation and washing procedures (permitting hybridisation between nucleotide sequences with 10% or less mismatches) were as described by Holden *et al*, (17). For non-stringent conditions (permitting hybridisation between sequences with 50% mismatches) filters were
30 hybridised overnight at 42°C in 10% formamide/0.25 M Na₂HPO₄/7% SDS

and the most stringent step was with 20 mM Na_2HPO_4 /1% SDS at 42°C. DNA fragments used as probes were labelled with [^{32}P]dCTP using the 'Radprime' system (Gibco-BRL) or with [digoxigenin-11]dUTP and detected using the Digoxigenin system (Boehringer Mannheim) according to the manufacturers' instructions, except that hybridisation was performed in the same solution as that used for radioactively labelled probes. Genomic DNA was prepared for Southern hybridisation as described previously (13).

Molecular Cloning and Nucleotide Sequencing. Restriction endonucleases and T4 DNA ligase were obtained from Gibco-BRL. General molecular biology techniques were as described in Sambrook *et al*, (18). Nucleotide sequencing was performed by the dideoxy chain termination method (19) using a T7 sequencing kit (Pharmacia). Sequences were assembled with the MacVector 3.5 software or AssemblyLIGN packages. Nucleotide and derived amino acid sequences were compared with those in the European Molecular Biology Laboratory (EMBL) and SwissProt databases using the BLAST and FASTA programs of the GCG package from the University of Wisconsin (version 8) (20) on the network service at the Human Genome Mapping Project Resource Centre, Hinxton, UK.

Virulence Tests. Groups of five female BALB/c mice (20-25g) were inoculated orally (p.o.) or intraperitoneally (i.p.) with 10-fold dilutions of bacteria suspended in physiological saline. For preparation of the inoculum, bacteria were grown overnight at 37°C in LB broth with shaking (50 rpm) and then used to inoculate fresh medium for various lengths of time until an optical density (OD) at 560 nm of 0.4 to 0.6 had been reached. For cell densities of 5×10^8 colony forming units (cfu) per ml and above, cultures were concentrated by centrifugation and resuspended in saline. The concentration of cfu/ml was checked by plating a dilution series of the inoculum onto LB agar plates. Mice were inoculated i.p. with 0.2 ml

volumes and p.o. by gavage with the same volume of inoculum. The LD₅₀ values were calculated after 28 days by the method of Reed and Meunch (21).

5 Results

Localisation of Transposon Insertions. The generation of a bank of *Salmonella typhimurium* mini Tn5 transposon mutants and the screen used to identify 43 mutants with attenuated virulence have been described previously (5). Transposons and flanking DNA regions were cloned from exconjugants by selection for kanamycin resistance or by inverse PCR. Nucleotide sequences of 300-600 bp of DNA flanking the transposons were obtained for 33 mutants. Comparison of these sequences with those in the DNA and protein databases indicated that 14 mutants resulted from transposon insertions into previously known virulence genes, 7 arose from insertions into new genes with similarity to known genes of the enterobacteria and 12 resulted from insertions into sequences without similarity to entries in the DNA and protein databases (ref. 5, Example 1 and this Example).

20

Three lines of evidence suggested that 16 of 19 transposon insertions into new sequences were clustered in three regions of the genome, initially designated A, B and C. First, comparing nucleotide sequences from regions flanking transposon insertion points with each other and with those in the databases showed that some sequences overlapped with one another or had strong similarity to different regions of the same gene. Second, Southern analysis of genomic DNA digested with several restriction enzymes and probed with restriction fragments flanking transposon insertion points indicated that some transposon insertions were located on the same restriction fragments. Third, when the same DNA probes were hybridised

30

to plaques from a *S. typhimurium* λ DNA library, the probes from mutants which the previous two steps had suggested might be linked were found to hybridise to the same λ DNA clones. Thus two mutants (P9B7 and P12F5) were assigned to cluster A, five mutants (P2D6, P9B6, P11C3, P11D10 and
5 P11H10) to cluster B and nine mutants (P3F4, P4F8, P7A3, P7B8, P7G2, P8G12, P9G4, P10E11 and P11B9) to cluster C (Figure 8).

Hybridisation of DNA probes from these three clusters to lysates from a set of *S. typhimurium* strains harbouring locked-in Mud-P22 prophages (15, 16)
10 showed that the three loci were all located in the minute 30 to 31 region (edition VIII, ref. 22) (Figure 7), indicating that the three loci were closely linked or constituted one large virulence locus. To determine if any of the λ clones covering clusters A, B and C contained overlapping DNA inserts, DNA fragments from the terminal regions of each clone were used as
15 probes in Southern hybridisation analysis of the other λ clones. Hybridising DNA fragments showed that several λ clones overlap and that clusters A, B and C comprise one contiguous region (Figure 8). DNA fragments from the ends of this region were then used to probe the λ library to identify further clones containing inserts representing the adjacent regions. No λ
20 clones were identified that covered the extreme right hand terminus of the locus so this region was obtained by cloning a 6.5 kb *EcoRI/XbaI* fragment from a lysate of the Mud-P22 prophage strain TT15244 (16).

Restriction mapping and Southern hybridisation analysis were then used to
25 construct a physical map of this locus (Figure 8). To distinguish this locus from the well characterised *inv/spa* gene cluster at minute 63 (edition VIII, ref. 22) (8, 9, 23, 24, 25, 26), we refer to the latter as virulence gene cluster 1 (VGC1) and have termed the new virulence locus VGC2. Figure 2 shows the position of two portions of DNA whose nucleotide sequence
30 has been determined ("Sequence 1" and "Sequence 2"). The nucleotide

sequence is shown in Figures 11 and 12.

Mapping the boundaries of VGC2 on the *S. typhimurium* chromosome.

Nucleotide sequencing of λ clone 7 at the left hand side of VGC2 revealed
5 the presence of an open reading frame (ORF) whose deduced amino acid
sequence is over 90% identical to the derived product of a segment of the
ydhE⁺ gene of *E. coli* and sequencing of the 6.5 kb *EcoRI/XbaI* cloned
fragment on the right hand side of VGC2 revealed the presence of an ORF
whose predicted amino acid sequence is over 90% identical to pyruvate
10 kinase I of *E. coli* encoded by the *pykF* gene (27). On the *E. coli*
chromosome *ydhE* and *pykF* are located close to one another, at minute 37
to 38 (28). Eleven non-overlapping DNA fragments distributed along the
length of VGC2 were used as probes in non-stringent Southern hybridisation
analysis of *E. coli* and *S. typhimurium* genomic DNA. Hybridising DNA
15 fragments showed that a region of approximately 40 kb comprising VGC2
was absent from the *E. coli* genome and localised the boundaries of VGC2
to within 1 kb (Figure 9). Comparison of the location of the *XbaI* site close
to the right hand end of VGC2 (Figure 8) with a map of known *XbaI* sites
(29) at the minute 30 region of the chromosome (22) enables a map position
20 of 30.7 minutes to be deduced for VGC2.

Structure of VGC2. Nucleotide sequencing of portions of VGC2 has
revealed the presence of 19 ORFs (Figure 8). The G+C content of
approximately 26 kb of nucleotide sequence within VGC2 is 44.6%,
25 compared to 47% for VGC1 (9) and 51-53% estimated for the entire
Salmonella genome (30).

The complete deduced amino acid sequences of ORFs 1-11 are similar to
those of proteins of type III secretion systems (6, 7), which are known to
30 be required for the export of virulence determinants in a variety of bacterial

pathogens of plants and animals (7). The predicted proteins of ORFs 1 - 8 (Figure 8) are similar in organisation and sequence to the products of the *yscN-U* genes of *Yersinia pseudotuberculosis* (31), to *invC/spaS* of the *inv/spa* cluster in VGC1 of *Salmonella typhimurium* (8, 9) and to *spa47/spa40* of the *spa/mxi* cluster of *Shigella flexneri* (32, 33, 34, 35,). For example the predicted amino acid sequence of ORF 3 (Figure 8) is 50% identical to YscS of *Y. pseudotuberculosis* (31), 34% identical to Spa9 from *S. flexneri* (35) and 37% identical to SpaQ of VGC1 of *S. typhimurium* (9). The predicted protein product of ORF9 is closely related to the LcrD family of proteins with 43% identity to LcrD of *Y. enterocolitica* (36), 39% identity to MxiA of *S. flexneri* (32) and 40% identity to InvA of VGC1 (23). Partial nucleotide sequences for the remaining ORFs shown in Figure 8 indicate that the predicted protein from ORF10 is most similar to *Y. enterocolitica* YscJ (37) a lipoprotein located in the bacterial outer membrane, with ORF11 similar to *S. typhimurium* InvG, a member of the PulD family of translocases (38). ORF12 and ORF13 show significant similarity to the sensor and regulatory subunits respectively, from a variety of proteins comprising two component regulatory systems (39). There is ample coding capacity for further genes between ORFs 9 and 10, ORFs 10 and 11, and between ORF 19 and the right hand end of VGC2.

VGC2 is conserved among and is specific to the *Salmonellae*. A 2.2 kb *PstI/HindIII* fragment located at the centre of VGC2 (probe B, Figure 8) lacking sequence similarity to entries in the DNA and protein databases was used as a probe in Southern hybridisation analysis of genomic DNA from *Salmonella* serovars and other pathogenic bacteria (Figure 10A). DNA fragments hybridising under non-stringent conditions showed that VGC2 is present in *S. aberdeen*, *S. gallinarum*, *S. cubana*, *S. typhi* and is absent from EPEC, EHEC, *Y. pestis*, *S. flexneri*, *V. cholera* and *S. aureus*. Thus VGC2 is conserved among and is likely to be specific to the *Salmonellae*.

To determine if the organisation of the locus is conserved among the *Salmonella* serovars tested, stringent Southern hybridisations with genomic DNA digested with two further restriction enzymes were carried out. Hybridising DNA fragments showed that there is some heterogeneity in the arrangement of restriction sites between *S. typhimurium* LT2 and *S. gallinarum*, *S. cubana* and *S. typhi* (Figure 10B). Furthermore, *S. gallinarum* and *S. typhi* contain additional hybridising fragments to those present in the other *Salmonellae* examined, suggesting that regions of VGC2 have been duplicated in these species.

10

VGC2 is required for virulence in mice. Previous experiments showed that the LD₅₀ values for i.p. inoculation of transposon mutants P3F4, P7G2, P9B7 and P11C3 were at least 100-fold greater than the wild type strain (5). In order to clarify the importance of VGC2 in the process of infection, the p.o. and i.p. LD₅₀ values for mutants P3F4 and P9B7 were determined (Table 1). Both mutants showed a reduction in virulence of at least five orders of magnitude by either route of inoculation in comparison with the parental strain. This profound attenuation of virulence by both routes of inoculation demonstrates that VGC2 is required for events in the infective process after epithelial cell penetration in BALB/c mice.

20

Table 1. LD₅₀ values of *S. typhimurium* strains.

Strain	LD ₅₀ (cfu)	
	i.p.	p.o.
12023 wild type	4.2	6.2 x 10 ⁴
P3F4	1.5 x 10 ⁶	> 5 x 10 ⁹
P9B7	> 1.5 x 10 ⁶	> 5 x 10 ⁹

25

cfu. colony forming units

Discussion

A hitherto unknown virulence locus in *S. typhimurium* of approximately 40 kb located at minute 30.7 on the chromosome by mapping the insertion points of a group of signature-tagged transposon mutants with attenuated virulence has been identified (5). This locus is referred to as virulence gene cluster 2 (VGC2) to distinguish it from the *inv/spa* virulence genes at 63 minutes (edition VIII, ref. 22) which we suggest be renamed VGC1. VGC1 and VGC2 both encode components of type III secretion systems. However, these secretion systems are functionally distinct.

Of 19 mutants that arose from insertions into new genes (ref. 5 and this example) 16 mapped to the same region of the chromosome. It is possible that mini-Tn5 insertion occurs preferentially in VGC2. Alternatively, as the negative selection used to identify mutants with attenuated virulence (5) was very stringent (reflected by the high LD₅₀ values for VGC2 mutants) it is possible that, among the previously unknown genes, only mutations in those of VGC2 result in a degree of attenuation sufficient to be recovered in the screen. The failure of previous searches for *S. typhimurium* virulence determinants to identify VGC2 might stem from reliance on cell culture assays rather than a live animal model of infection. A previous study which identified regions of the *S. typhimurium* LT2 chromosome unique to *Salmonellae* (40) located one such region (RF333) to minutes 30.5 - 32. Therefore, RF333 may correspond to VGC2, although it was not known that RF333 was involved in virulence determination.

Comparisons with the type III secretion systems encoded by the virulence plasmids of *Yersinia* and *Shigella* as well as with VGC1 of *Salmonella* indicates that VGC2 encodes the basic structural components of the secretory apparatus. Furthermore, the order of ORFs 1-8 in VGC2 is the

- same as the gene order in homologues in *Yersinia*, *Shigella* and VGC1 of *S. typhimurium*. The fact that the organisation and structure of the VGC2 secretion system is no more closely related to VGC1 than to the corresponding genes of *Yersinia*, together with the low G+C content of VGC2 suggests that VGC2, like VGC1 (40, 41, 42) was acquired independently by *S. typhimurium* via horizontal transmission. The proteins encoded by ORFs 12 and 13 show strong similarity to bacterial two component regulators (39) and could regulate either ORFs 1-11 and/or the secreted proteins of this system.
- Many genes in VGC1 have been shown to be important for entry of *S. typhimurium* into epithelial cells. This process requires bacterial contact (2) and results in cytoskeletal rearrangements leading to localised membrane ruffling (43, 44). The role of VGC1 and its restriction to this stage of the infection is reflected in the approximately 50-fold attenuation of virulence in BALB/c mice inoculated p.o. with VGC1 mutants and by the fact that VGC1 mutants show no loss of virulence when administered i.p. (8). The second observation also explains why no VGC1 mutants were obtained in our screen (5). In contrast, mutants in VGC2 are profoundly attenuated following both p.o. and i.p. inoculation. This shows that, unlike VGC1, VGC2 is required for virulence in mice after epithelial cell penetration, but these findings do not exclude a role for VGC1 in this early stage of infection.

- Thus in summary mapping the insertion points of 16 signature-tagged transposon mutants on the *Salmonella typhimurium* chromosome led to the identification of a 40 kb virulence gene cluster at minute 30.7. This locus is conserved among all other *Salmonella* species examined, but not present in a variety of other pathogenic bacteria or in *Escherichia coli* K12. Nucleotide sequencing of a portion of this locus revealed 11 open reading frames whose predicted proteins encode components of a type III secretion

system. To distinguish between this and the type III secretion system encoded by the *inv/spa* invasion locus we refer to the *inv/spa* locus as virulence gene cluster 1 (VGC1) and the new locus as VGC2. VGC2 has a lower G+C content than that of the *Salmonella* genome and is flanked by genes whose products share greater than 90% identity with those of the *E. coli ydhE* and *pykF* genes. Thus VGC2 was probably acquired horizontally by insertion into a region corresponding to that between the *ydhE* and *pykF* genes of *E. coli*. Virulence studies of VGC2 mutants have shown them to be attenuated by at least five orders of magnitude compared with the wild type strain following oral or intraperitoneal inoculation.

References for this Example

1. Carter, P.B. & Collins, F.M. (1974) *J. Exp. Med.* **139**, 1189-1203.
- 15 2. Takeuchi, A. (1967) *Am. J. Pathol.* **50**, 109-136.
3. Finlay, B.B. (1994) *Curr. Top. Microbiol. Immunol.* **192**, 163-185.
4. Groisman, E.A. & Ochman, H. (1994) *Trends Microbiol.* **2**, 289-293.
5. Hensel, M., Shea, J.E., Gleeson, C., Jones, M.D., Dalton, E. &
20 Holden, D.W. (1995) *Science* **269**, 400-403.
6. Salmond, G.P.C. & Reeves, P.J. (1993) *Trends Biochem. Sci.* **18**, 7-12.
7. Van Gijsegem, F., Genin, S. & Boucher, C. (1993) *Trends Microbiol.* **1**, 175-180.
- 25 8. Galan, J.E. & Curtiss, R. (1989) *Proc. Natl. Acad. Sci. U.S.A.* **86**, 6383-6387.
9. Groisman, E.A. & Ochman, H. (1993) *EMBO J.* **12**, 3779-3787.
10. Ginocchio, C.C., Olmsted, S.B., Wells, C.L. & Galan, J.E. (1994) *Cell* **76**, 717-724.
- 30 11. de Lorenzo, V. & Timmis, K.N. (1994) *Methods Enzymol.* **264**.

- 386-405.
12. Davis, R.H., Botstein, D. & Roth, J.R. (1980) *Advanced Bacterial Genetics*, Cold Spring Harbor Laboratory, Cold Spring Harbor, N.Y.
 - 5 13. Ausubel, F.M., Brent, R., Kingston, R.E., Moore, D.D., Seidman, J.G., Smith, J.A. and Struhl, K. (1987) *Current Protocols in Molecular Biology* Vol 4 John Wiley and Sons, Inc, New York.
 14. Maurer, R., Osmond, B.C., Shekhtman, E., Wong, A. & Botstein, D. (1984) *Genetics* **108**, 1-23.
 - 10 15. Youderain, P., Sugiono, P., Brewer, K.L., Higgins, N.P. & Elliott, T. (1988) *Genetics* **118**, 581-592.
 16. Benson, N.R. & Goldman, B.S. (1992) *J. Bacteriol.* **174**, 1673-1681.
 17. Holden, D.W., Kronstad, J.W. & Leong, S. (1989) *EMBO J.* **8**, 1927-1934.
 - 15 18. Sambrook, J., Fritsch, E.F. & Maniatis, T. (1989) *Molecular cloning: a laboratory manual* (Cold Spring Harbor Laboratory, Cold Spring Harbor, New York).
 19. Sanger, F., Nicklen, S. & Coulson, A.R. (1977) *Proc. Natl. Acad. Sci. U.S.A.* **74**, 5463-5467.
 - 20 20. Devereux, J., Hearberli, P. & Smithies, O. (1984) *Nucl. Acids Res.* **12**, 387-399.
 21. Reed, L.J. & Muench, H. (1938) *Am. J. Hyg.* **27**, 493-497.
 22. Sanderson, K.E., Hessel, A. & Rudd, K.E. (1995) *Microbiol. Rev.* **59**, 241-303.
 - 25 23. Galan, J.E., Ginocchio, C. & Costeas, P. (1992) *J. Bacteriol.* **174**, 4338-4349.
 24. Ginocchio, C., Pace, J. & Galan, J.E. (1992) *Proc. Natl. Acad. Sci. U.S.A.* **89**, 5976-5980.

25. Eichelberg, K., Ginocchio, C.C. & Galan, J.E. (1994) *J. Bacteriol.* **176**, 4501-4510.
26. Collazo, C.M., Zierler, M.K. & Galan, J.E. (1995) *Mol. Microbiol.* **15**, 25-38.
- 5 27. Ohara, O., Dorit, R.L. & Gilbert, W. (1989) *Proc. Natl. Acad. Sci. USA* **86**, 6883-6887.
28. Bachman, B. (1990) *Micro. Rev.* **54**, 130-197.
29. Liu, S.L., Hessel, A. & Sanderson, K.E. (1993) *J. Bacteriol.* **175**, 4104-4120.
- 10 30. Fasman, G.D. (1976) *CRC Handbook of Biochemistry and Molecular Biology*, CRC Press, Cleveland.
31. Bergman, T., Erickson, K., Galyov, E., Persson, C. & Wolf Watz, H. (1994) *J. Bacteriol.* **176**, 2619-2626.
32. Andrews, G.P. & Maurelli, A.T. (1992) *Infect. Immun.* **60**, 3287-3295.
- 15 33. Allaoui, A., Sansonetti, P.J. & Parsot, C. (1993) *Mol. Microbiol.* **7**, 59-68.
34. Venkatesan, M., Buysse, J.M. & Oaks, E.V. (1992) *J. Bacteriol.* **174**, 1990-2001.
- 20 35. Sasakawa, C., Komatsu, K., Tobe, T., Suzuki, T. & Yoshikawa, M. (1993) *J. Bacteriol.* **175**, 2334-2346.
36. Plano, G.V., Barve, S.S. & Straley, S.C. (1991) *J. Bacteriol.* **173**, 7293-7303.
37. Michiels, T., Vanooteghem, J.C., Lambert de Rouvroit, C., China, B., Gustin, A., Boudry, P. & Cornelis, G.R. (1991) *J. Bacteriol.* **173**, 4994-5009.
- 25 38. Kaniga, K., Bossio, J.C. & Galan, J.E. (1994) *Mol. Microbiol.* **13**, 555-568.
39. Ronson, C.W., Nixon, B.T. & Ausubel, F.M. (1987) *Cell* **49**, 579-581.
- 30

40. Groisman, E.A., Sturmoski, M.A., Solomon, F.R., Lin, R. & Ochman, H. (1993) *Proc. Natl. Acad. Sci. U.S.A.* **90**, 1033-1037.
41. Altmeyer, R.M., McNern, J.K., Bossio, J.C., Rosenshine, I., Finlay, B.B. & Galan, J.E. (1993) *Mol. Microbiol.* **7**, 89-98.
- 5 42. Li, J., Ochman, H., Groisman, E.A., Boyd, E.F., Solomon, F., Nelson, K. & Selander, R.K. (1995) *Proc. Natl. Acad. Sci. USA* **92**, 7252-7256.
43. Finlay, B.B. & Rauschkowski, S. (1991) *J. Cell Sci.* **99**, 283-296.
44. Francis, C.L., Starnbach, M.N. & Falkow, S. (1992) *Mol.*
10 *Microbiol.* **6**, 3077-3087.

Example 5: Identification of virulence genes in *Streptococcus pneumoniae*

(a) Mutagenesis

- 15 In the absence of a convenient transposon system, the most efficient way of creating tagged mutants of *Streptococcus pneumoniae* is to use insertion-duplication mutagenesis (Morrison *et al* (1984) *J. Bacteriol.* **159**, 870). Random *S. pneumoniae* DNA fragments of 200-400 bp will be generated by genomic DNA digestion with a restriction enzyme or by
- 20 physical shearing by sonication followed by gel fractionation and DNA end-repair using T4 DNA polymerase. The fragments are ligated into plasmid pJDC9 (Pearce *et al* (1993) *Mol. Microbiol.* **9**, 1037 which carries the *erm* gene for erythromycin selection in *E. coli* and *S. pneumoniae*), previously modified by incorporation of DNA sequence tags into one of the
- 25 polylinker cloning sites. The size of cloned *S. pneumoniae* DNA is sufficient to ensure homologous recombination, and reduces the possibility of generating an unrepresentative library in *E. coli* (expression of *S. pneumoniae* proteins can be toxic to *E. coli*). Alternative vectors carrying different selectable markers are available and can be used in place of
- 30 pJDC9. Tagged plasmids carrying DNA fragments are introduced to an

appropriate *S. pneumoniae* strain selected on the basis of serotype and virulence in a murine model of pneumococcal pneumonia. Regulation of competence for genetic transformation in *S. pneumoniae* is governed by competence factor, a peptide of 17 amino acids which has been characterized recently by Don Morrison's group at the University of Illinois at Chicago and which is described Havarstein, Coomaraswamy and Morrison (1995) *Proc. Natl. Acad. Sci. USA* 92, 11140-11144. Incorporation of minute quantities of this peptide in transformation experiments leads to very efficient transformation frequencies in some encapsulated clinical isolates of *S. pneumoniae*. This overcomes a major hurdle in pneumococcal molecular genetics and the availability of the peptide greatly facilitates the construction of *S. pneumoniae* mutant banks and allows flexibility in choosing the strain(s) to be mutated. A proportion of transformants are analysed to verify homologous integration of the plasmid sequences, and checked for stability. The very low level of reversion associated with mutants generated by insertion-duplication is minimized by the fact that the duplicated regions will be short (200-400 bp); however if the level of reversion is unacceptably high, antibiotic selection is maintained during growth of the transformants in culture and during growth in the animal.

(b) Animal model

The *S. pneumoniae* mutant bank is organized into pools for inoculation into Swiss and/or C57B1/6 mice. Preliminary experiments are conducted to determine the optimum complexity of the pools and the optimum inoculum level. One attractive model utilises inocula of 10^5 cfu, delivered by mouth to the trachea (Veber *et al* (1993) *J. Antimicrobial Chemotherapy* 32, 473). Swiss mice develop acute pneumonia within 3-4 days, and C57B1/6 mice develop subacute pneumonia within 8-10 days. These pulmonary models of infection yield 10^8 cfu/lung (Veber *et al* (1993) *J. Antimicrobial*

Chemotherapy 32, 473) at the time of death. If required, mice are also injected intraperitoneally for the identification of genes required for bloodstream infection (Sullivan *et al* (1993) *Antimicrobial Agents and Chemotherapy* 37, 234).

5

(c) Virulence gene identification

Once the parameters of the infection model are optimized, a mutant bank consisting of several thousand strains is subjected to virulence tests. Mutants with attenuated virulence are identified by hybridisation analysis, using labelled tags from the 'input' and 'recovered' pools as probes. If *S. pneumoniae* DNA cannot be colony blotted easily, chromosomal DNA is liberated chemically or enzymatically in the wells of microtitre dishes prior to transfer onto nylon membranes using a dot-blot apparatus. DNA flanking the integrated plasmid is cloned by plasmid rescue in *E. coli* (Morrison *et al* (1984) *J. Bacteriol.* 159, 870), and sequenced. Genomic DNA libraries are constructed in appropriate vectors maintained in either *E. coli* or a Gram-positive host strain, and are probed with restriction fragments flanking the integrated plasmid to isolate cloned virulence genes which is then fully sequenced and subjected to detailed functional analysis.

20

Example 6: Identification of virulence genes in *Enterococcus faecalis*

(a) Mutagenesis

Mutagenesis of *E. faecalis* is accomplished using plasmid pAT112 or a derivative, developed for this purpose. pAT112 carries genes for selection in both Gram-negative and Gram-positive bacteria, and the *att* site of Tn/545. It therefore requires the presence in the host strain of the integrase for transposition, and stable, single copy insertions are obtained if the host does not contain an excisionase gene (Trieu-Cuot *et al* (1991) *Gene* 106, 21). Recovery of DNA flanking the integrated plasmid is accomplished by

30

restriction digestion of genomic DNA, intramolecular ligation and transformation of *E. coli*. The presence of single sites for restriction enzymes in pAT112 and its derivatives will (Trieu-Cuot *et al* (1991) *Gene* 106, 21) allows the incorporation of DNA sequence tags prior to transfer
5 to a virulent strain of *E. faecalis* carrying plasmid pAT145 (to provide the integrase function) by either conjugation, electroporation or transformation (Trieu-Cuot *et al* (1991) *Gene* 106, 21; Wirth *et al* (1986) *J. Bacteriol.* 165, 831).

10 (b) Animal model

A large number of insertion mutants are analysed for random integration of the plasmid by isolating DNA from transciipients, restriction enzyme digestion and Southern hybridisation. Individual mutants are stored in the wells of microtitre dishes, and complexity and size of pooled inocula are
15 optimised prior to screening of the mutant bank. Two different models of infection caused by *E. faecalis* are employed. The first is a well established rat model of endocarditis, involving tail vein injection of up to 10^8 cfu of *E. faecalis* into animals that have a catheter inserted across the aortic valve (Whitman *et al* (1993) *Antimicrobial Agents and Chemotherapy* 37, 1069).
20 Animals are sacrificed at various times after inoculation, and bacterial vegetations on the aortic valve are excised, homogenized and plated to culture medium to recover bacterial colonies. Virulent bacteria are also recovered from the blood at various times after inoculation. The second model is of peritonitis in mice, following intraperitoneal injection of up to
25 10^9 cfu of *E. faecalis* (Chenoweth *et al* (1990) *Antimicrobial Agents and Chemotherapy* 34, 1800). As with the *S. pneumoniae* model, preliminary experiments are done to establish the optimum complexity of the pools and the optimum inoculum level, prior to screening the mutant bank.

(c) Virulence gene identification

Isolation of DNA flanking the site of integration of pAT112 using its *E. coli* origin of replication is simplified by the lack of sites for most of the commonly used 6 bp recognition restriction enzymes in the vector.

- 5 Therefore DNA from the strains of interest are digested with one of these enzymes, self-ligated, transformed into *E. coli* and sequenced using primers based on the sequences adjacent to the *att* sites on the plasmid. A genomic DNA library of *E. faecalis* are probed with sequences of interest to identify intact copies of virulence genes which are then sequenced.

10

Example 7: Identification of virulence genes in *Pseudomonas aeruginosa*

(a) Mutagenesis

- 15 Since transposon Tn5 has been used by others to mutagenise *Pseudomonas aeruginosa*, and the mini-Tn5 derivative that was used for the identification of *Salmonella typhimurium* virulence genes (Example 1) is reported to have broad utilisation among Gram-negative bacteria, including several pseudomonads (DeLorenzo and Timaris (1994) *Methods Enzymol.* 264,
20 386), a *P. aeruginosa* mutant bank is constructed using our existing pool of signature tagged mini-Tn5 transposons by conjugal transfer of the suicide vector to one or more virulent (and possibly mucoid) recipient strains. This approach represents a significant time saving. Other derivatives of Tn5 designed specifically for *P. aeruginosa* mutagenesis (Rella *et al* (1985) *Gene*
25 33, 293), may alternatively be employed with the mini Tn5 transposon.

(b) Animal model and virulence gene identification

- The bank of *P. aeruginosa* insertion mutants is screened for attenuated virulence in a chronic pulmonary infection model in rats. Suspensions of
30 *P. aeruginosa* cells are introduced into a bronchus following tracheotomy,

and disease develops over a 30 day period (Woods *et al* (1982) *Infect. Immun.* 36, 1223). Bacteria are recovered by plating lung homogenates to laboratory medium and sequence tags from these are used to probe DNA colony blots of bacteria used as the inoculum. It is also possible to subject
5 the mutant bank to virulence tests in a model of endogenous bacteremia (Hirakata *et al* (1992) *Antimicrobial Agents and Chemotherapy* 36, 1198), and cystic fibrosis (Davidson *et al* (1995) *Nature Genetics* 9, 351) in mice. Cloning and sequencing of DNA flanking the transposons is done as described in Example 1. Genomic DNA libraries for the isolation and
10 sequencing of intact copies of the genes are constructed in the laboratory by standard methods.

Example 8: Identification of virulence genes in *Aspergillus fumigatus*

15 (a) Mutagenesis

The functional equivalent of transposon mutagenesis in fungi is restriction enzyme mediated integration (REMI) of transforming DNA (Schiestl and Petes (1991) *Proc. Natl. Acad. Sci.* 88, 7585). In this process, fungal cells are transformed with DNA fragments carrying a selectable marker in the
20 presence of a restriction enzyme, and single copy integrations occur at different genomic sites, defined by the target sequence of the restriction enzyme. REMI has already been used successfully to isolate virulence genes of *Cochliobolus* (Lu *et al* (1994) *Proc. Natl. Acad. Sci. USA* 91, 12649) and *Ustilago* (Bolker *et al* (1995) *Mol. Gen. Genet.* 248, 547), and
25 have shown that incorporation of active restriction enzyme with a plasmid encoding hygromycin resistance leads to single and apparently random integration of the linear plasmid into the *A. fumigatus* genome. Sequence tags are introduced into a convenient site in one of two vectors for hygromycin resistance, and used to transform a clinical isolate of *A.*
30 *fumigatus*.

(b) Animal model and virulence gene identification

The low-dose model of aspergillosis in neutropenic mice in particular closely matches the course of pulmonary disease in humans (Smith *et al* (1994) *Infect. Immun.* 62, 5247). Mice are inoculated intranasally with up to 1,000,000 conidiospores/mouse, and virulent fungal mutants are recovered 7-10 days later by using lung homogenates to inoculate liquid medium. Hyphae are collected after a few hours, from which DNA is extracted for amplification and labelling of tags to probe colony blots of DNA from the pool of transformants comprising the inoculum. DNA from the regions flanking the REMI insertion points are cloned by digesting the transformant DNA with a restriction enzyme that cuts outside the REMI vector, self ligation and transformation of *E. coli*. Primers based on the known sequence of the plasmid are used to determine the adjacent *A. fumigatus* DNA sequences. To prove that the insertion of the vector was the cause of the avirulent phenotype, the recovered plasmid is recut with the same restriction enzyme used for cloning, and transformed back into the wild-type *A. fumigatus* parent strain. Transformants that have arisen by homologous recombination are then subjected to virulence tests.

REFERENCES (other than for Example 4)

- Ausubel, F.M., Brent, R., Kingston, R.E., Moore, D.D., Seidman, J.G., Smith, J.A. and Struhl, K. (1987) *Current Protocols in Molecular Biology*, New York: John Wiley and Sons.
- Buchmeier, N.A., Lipps, C.J., So, M.Y. and Heffron, F. (1993) Recombination-deficient mutants of *Salmonella typhimurium* are avirulent and sensitive to the oxidative burst of macrophages. *Mol. Microbiol.* 7, 933-936.
- 10 Carter, P.B. and Collins, F.M. (1974) The route of enteric infection in normal mice. *J. Exp. Med.* 139, 1189-1203.
- de Lorenzo, V. and Timmis, K.N. (1994) Analysis and construction of stable phenotypes in Gram-negative bacteria with Tn5-and Tn10-derived minitransposons. *Methods Enzymol.* 264, 386-405.
- 15 de Lorenzo, V., Herrero, M., Jakubzik, U. and Timmis, K.N. (1990) Mini-Tn5 transposon derivatives for insertion mutagenesis, promoter probing, and chromosomal insertion of cloned DNA in gram-negative eubacteria. *J. Bacteriol.* 172, 6568-6572.
- 20 Fields, P.I., Groisman, E.A. and Heffron, F. (1989) A *Salmonella* locus that controls resistance to microbicidal proteins from phagocytic cells. *Science* 243, 1059-1062.
- Finlay, B.B., Starnbach, M.N., Francis, C.L., Stocker, B.A., Chatfield, S., Dougan, G. and Falkow, S. (1988) Identification and characterization of Tn ϕ oA mutants of *Salmonella* that are unable to pass through a

polarized MDCK epithelial cell monolayer. *Mol. Microbiol.* 2, 757-766.

Groisman, E.A., Chiao, E., Lipps, C.J., Heffron, F. (1989) *Salmonella typhimurium phoP* virulence gene is a transcriptional regulator. *Proc. Natl. Acad. Sci. USA.* 86, 7077-7081.

- 5 Groisman, E.A. and Ochman, H. (1994) How to become a pathogen. *Trends Microbiol.* 2, 289-293.

Groisman, E.A. and Saier, M.H., Jr. (1990) *Salmonella* virulence: new clues to intramacrophage survival. *Trends Biochem. Sci.* 15, 30-33.

- 10 Herrero, M., de Lorenzo, V. and Timmis, K.N. (1990) Transposon vectors containing non-antibiotic resistance selection markers for cloning and stable chromosomal insertion of foreign genes in Gram-negative bacteria. *J. Bacteriol.* 172, 6557-6567.

Holden D.W., Kronstad J.W., Leong S.A. (1989) Mutation in a heat-regulated hsp70 gene of *Ustilago maydis*. *EMBO J.* 8, 1927-1934.

- 15 Holland J., Towner K.J., Williams P. (1992) Tn916 insertion mutagenesis in *Escherichia coli* and *Haemophilus influenzae* type b following conjugative transfer. *J. Gen. Microbiol.* 138, 509-515.

- Mahan, M.J., Slauch, J.M., Mekalanos, J.J. (1993) Selection of bacterial virulence genes that are specifically induced in host tissues. *Science* 259, 686-688.

Miller, S.I., Kukral, A.M. and Mekalanos, J.J. (1989a) A two-component regulatory system (*phoP phoQ*) controls *Salmonella typhimurium* virulence. *Proc. Natl. Acad. Sci. USA.* 86, 5054-5058.

Miller, I., Maskell, D., Hormaeche, C., Johnson, K., Pickard, D. and Dougan, G. (1989b) Isolation of orally attenuated *Salmonella typhimurium* following *TnphoA* mutagenesis. *Infect. Immun.* **57**, 2758-2763.

- 5 Miller, V.L. and Mekalanos, J.J. (1988) A novel suicide vector and its use in construction of inversion mutations: osmoregulation of outer membrane proteins and virulence determinants in *Vibrio cholerae* requires *toxR*. *J. Bacteriol.* **170**, 2575-2583.

- 10 Sambrook, J., Fritsch, E.F. and Maniatis, T. (1989) *Molecular cloning: a laboratory manual*, Cold Spring Harbor, New York: Cold Spring Harbor Laboratory.

Sanger, F., Nicklen, S. and Coulson, A.R. (1977) DNA sequencing with chain terminating inhibitors. *Proc. Natl. Acad. Sci. USA.* **74**, 5463-5467.

- 15 Schiestl R.H., and Petes T.D. (1991) Integration of DNA fragments by illegitimate recombination in *Saccharomyces cerevisiae*. *Proc. Natl. Acad. Sci USA.* **88**, 7585-7589.

CLAIMS

1. A method for identifying a microorganism having a reduced adaptation to a particular environment comprising the steps of:

- 5 (1) providing a plurality of microorganisms each of which is independently mutated by the insertional inactivation of a gene with a nucleic acid comprising a unique marker sequence so that each mutant contains a different marker sequence, or clones of the said microorganism;
- (2) providing individually a stored sample of each mutant
10 produced by step (1) and providing individually stored nucleic acid comprising the unique marker sequence from each individual mutant;
- (3) introducing a plurality of mutants produced by step (1) into the said particular environment and allowing those microorganisms which are able to do so to grow in the said environment;
- 15 (4) retrieving microorganisms from the said environment or a selected part thereof and isolating the nucleic acid from the retrieved microorganisms;
- (5) comparing any marker sequences in the nucleic acid isolated in step (4) to the unique marker sequence of each individual mutant stored
20 as in step (2); and
- (6) selecting an individual mutant which does not contain any of the marker sequences as isolated in step (4).

2. A method according to Claim 1 wherein the plurality of
25 microorganisms as defined in step (1) is produced from a plurality of microorganisms, each of which comprises a nucleic acid comprising a unique marker sequence, by changing their condition from a first given condition to a second given condition wherein (a) in the first given condition the said nucleic acid comprising a unique marker is maintained episomally
30 and (b) in the second given condition the said nucleic acid comprising a

unique marker sequence insertionally inactivates a gene.

3. A method according to Claims 1 or 2 further comprising the steps:
 - (1A) removing auxotrophs from the plurality of mutants produced
- 5 in step (1); or
 - (6A) determining whether the mutant selected in step (6) is an auxotroph; or
 - both (1A) and (6A).
- 10 4. A method of identifying a gene which allows a microorganism to adapt to a particular environment, the method comprising the method of any one of Claims 1 to 3 followed by the step:
 - (7) isolating the insertionally-inactivated gene from the individual mutant selected in step (6).
- 15 5. A method according to Claim 4 further comprising the step:
 - (8) isolating from a wild-type microorganism the corresponding wild-type gene using the insertionally-inactivated gene isolated in step (7) as a probe.
- 20 6. A method according to any one of Claims 1 to 5 wherein the particular environment is a differentiated multicellular organism.
7. A method according to Claim 6 wherein the multicellular organism
- 25 is a plant.
8. A method according to Claim 6 wherein the multicellular organism is a non-human animal.
- 30 9. A method according to Claim 8 wherein the animal is a mouse, rat.

rabbit, dog or monkey.

10. A method according to Claim 9 wherein the animal is a mouse.
- 5 11. A method according to any one of Claims 6 to 10 wherein in step (4) the microorganisms are retrieved from the said environment at a site remote from the site of introduction in step (3).
12. A method according to any one of Claims 8 to 10 wherein in step (3)
10 the microorganism is introduced orally or intraperitoneally.
13. A method according to Claim 12 when dependent on Claims 8 or 9 wherein in step (4) the microorganisms are retrieved from the spleen.
- 15 14. A method according to any one of the preceding claims wherein the microorganism is a bacterium.
15. A method according to any one of Claims 1 to 13 wherein the
20 microorganism is a fungus.
16. A method according to Claim 7 wherein the microorganism is a bacterium pathogenic to plants.
17. A method according to Claim 7 wherein the microorganism is a
25 fungus pathogenic to plants.
18. A method according to any one of Claims 8 to 10 wherein the microorganism is a bacterium pathogenic to animals.
- 30 19. A method according to any one of Claims 8 to 10 wherein the

microorganism is a fungus pathogenic to animals.

20. A method according to Claim 18 wherein the bacterium is any one of *Bordetella pertussis*, *Campylobacter jejuni*, *Clostridium botulinum*,
5 *Escherichia coli*, *Haemophilus ducreyi*, *Haemophilus influenzae*,
Helicobacter pylori, *Klebsiella pneumoniae*, *Legionella pneumophila*,
Listeria spp., *Neisseria gonorrhoeae*, *Neisseria meningitidis*, *Pseudomonas*
spp., *Salmonella* spp., *Shigella* spp., *Staphylococcus aureus*, *Streptococcus*
pyogenes, *Streptococcus pneumoniae*, *Vibrio* spp., and *Yersinia pestis*.

10

21. A method according to Claim 19 wherein the fungus is any one of
Aspergillus spp., *Cryptococcus neoformans* and *Histoplasma capsulatum*.

22. A method according to any one of the preceding claims wherein in
15 step (1) the gene is insertionally inactivated using a transposon or
transposon like element or other DNA sequence carrying a unique marker
sequence.

23. A method according to any one of the preceding claims wherein in
20 step (1) each different marker sequence is flanked on either side by
sequences common to each said nucleic acid.

24. A method according to Claim 23 wherein in step (2) the nucleic acid
comprising the unique marker is isolated using DNA amplification
25 techniques and oligonucleotide primers which hybridise to the said common
sequences.

25. A method according to Claim 23 or 24 wherein in step (4) the
nucleic acid comprising a plurality of said marker sequences is isolated
30 using DNA amplification techniques and oligonucleotide primers which

hybridise to the said common sequences.

26. A microorganism obtained using the method of any one of the preceding claims.

5

27. A microorganism comprising a mutation in a gene identified using the method of Claim 5.

28. A microorganism obtained according to Claim 26, when dependent
10 on Claim 8, or Claim 27 for use in a vaccine.

29. A vaccine comprising a microorganism according to Claim 26, when dependent on Claim 8, or Claim 27 and a pharmaceutically-acceptable carrier.

15

30. A gene obtained using the method of Claims 4 or 5.

31. A gene according to Claim 30 which is isolated from the *Salmonella typhimurium* genome and hybridises to the sequence shown in Figure 5
20 under stringent conditions.

32. A gene according to Claim 30 which is isolated from the *Salmonella typhimurium* genome and hybridises to a sequence shown in Figure 6 under stringent conditions.

25

33. A polypeptide encoded by a gene according to any one of Claims 30 to 32.

34. A method of identifying a compound which reduces the ability of a
30 microorganism to adapt to a particular environment comprising the step of

selecting a compound which interferes with the function of a gene according to any one of Claims 30 to 32 or a polypeptide according to Claim 33.

35. A compound identifiable by the method of Claim 34.

5

36. A compound according to Claim 35 wherein the particular environment is a host organism.

37. A compound according to Claim 36 wherein the host organism is a
10 plant.

38. A compound according to Claim 36 wherein the host organism is an animal.

15 39. Use of a compound according to any one of Claim 36 to Claim 38 for treating infection of said host organism with said microorganism.

40. A molecule which selectively interacts with, and substantially inhibits the function of, a gene according to any one of Claims 30 to 32 or a nucleic
20 acid product thereof.

41. A molecule according to Claim 40 which is an antisense nucleic acid or nucleic acid derivative.

25 42. A molecule according to Claim 40 or 41 which is an antisense oligonucleotide.

43. A molecule according to any one of Claims 40 to 42 for use in medicine.

30

44. A method of treating a host which has, or is susceptible to, an infection with a microorganism, the method comprising administering an effective amount of a molecule or compound according to Claim 36 or 40 wherein said gene is present in said microorganism, or a close relative of said microorganism.

45. A pharmaceutical composition comprising a molecule or compound according to Claim 38 or 40 and a pharmaceutically acceptable carrier.

46. The VGC2 DNA of *Salmonella typhimurium* or a part thereof, or a variant of said DNA or a variant of a part thereof.

47. A mutant bacterium wherein if the bacterium normally contains a gene that is the same as or equivalent to a gene in VGC2, said gene is mutated or absent in said mutant bacterium.

48. A method of making a bacterium according to Claim 47.

49. Use of a mutant bacterium according to Claim 47 in a vaccine.

50. A pharmaceutical composition comprising a bacterium according to Claim 47 and a pharmaceutically acceptable carrier.

51. A polypeptide encoded by VGC2 DNA of *Salmonella typhimurium* or a part thereof, or a variant of said polypeptide or a variant of a part thereof.

52. A method of identifying a compound which reduces the ability of a bacterium to infect or cause disease in a host comprising the step of selecting a compound which interferes with the function of a gene in VGC2

according to Claim 46 or a polypeptide according to Claim 51.

53. A compound identifiable by the method of Claim 52.

5 54. A molecule which selectively interacts with, and substantially inhibits the function of, a gene in VGC2 of *Salmonella typhimurium* or a nucleic product thereof.

10 55. A molecule or compound according to Claim 53 or 54 for use in medicine.

56. Any novel feature or combination of features disclosed herein.

DNA sequence tag

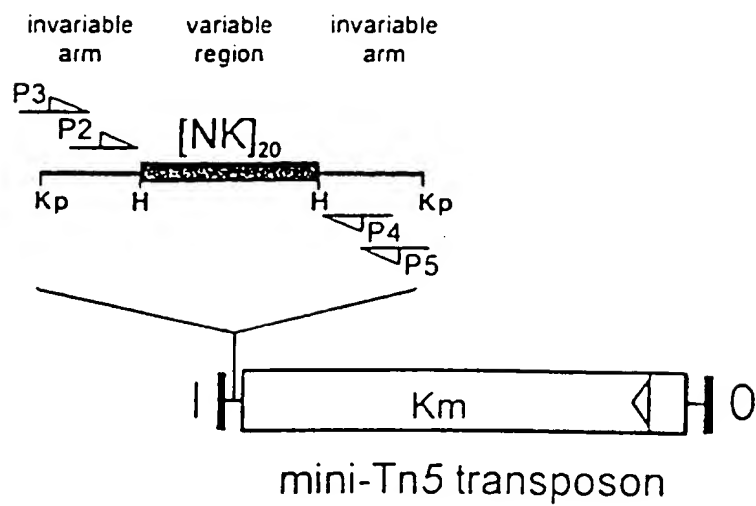


Figure 1a

2/39

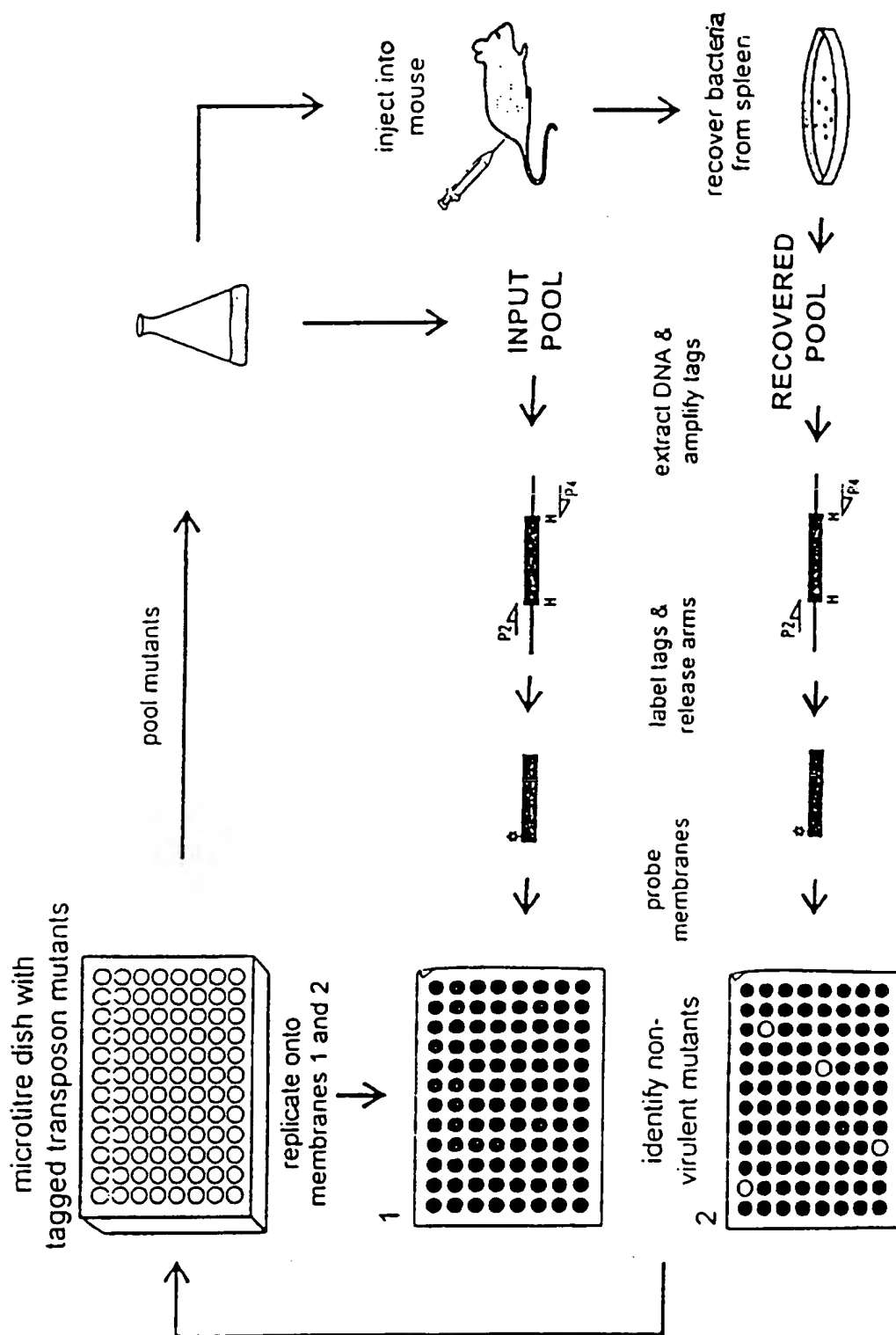


Figure 1b

3/39

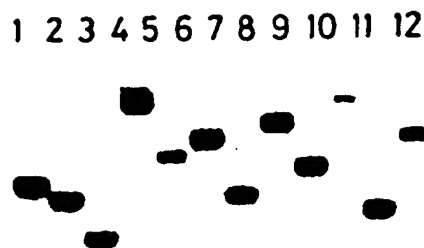


Fig. 2

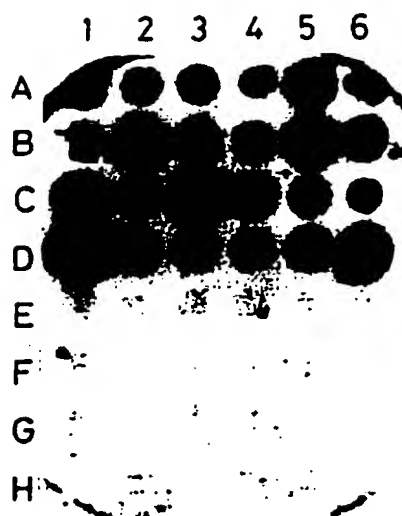


Fig. 3

4/39

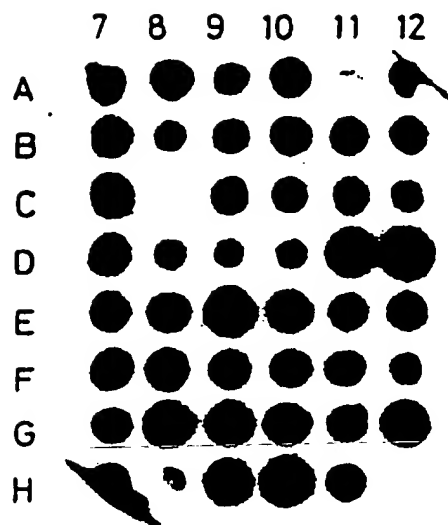
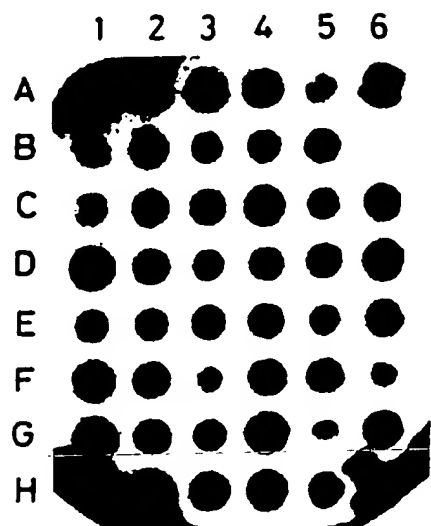
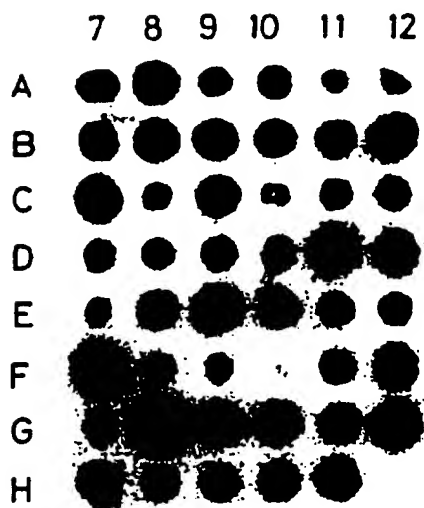
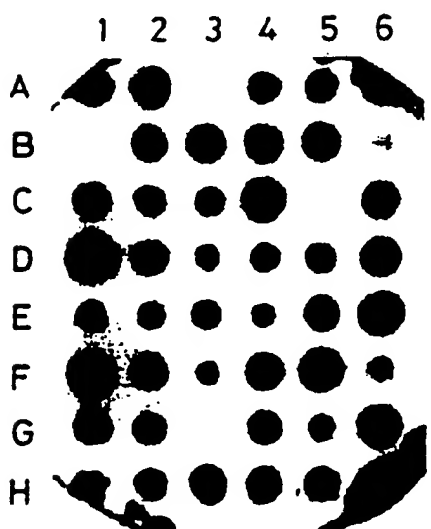
Inoculum pattern**Spleen pattern**

Fig. 4

Name. mpcc2 1

Minus Strand HSPs:

Minus Strand HSPs:

Fetch → Gb_ba:Ecoclp
- OK then type J Biol Chem 265, 12536.
(1990)

Figure 5

6/39

A) n w virulence factors with similarity to sequenced genes:

1. p1F10

similarity to *clpP* (*E.coli*)

(Figure 5 of application)

2. p2D6

similarity to *lcrD* (*Yersinia spp.*)

sequence p2D6_1_I

GGTCTTAATGTACGGGCATGGTCTGCATCGATAACTCCGGCAGCAAATCGCCATCGATACTCATTGT
 TTGGCTGGCATCCCATCAAGCGAGAAACGTGCGCTAACTTCCGCCACCCTCTCGATACCTTTTGTAAATG
 ACAATAAATTCACGATAGTAATGATGGTAAATACGACCAACCAACGGTGAGATTTCCTCCTACGACA
 AACTTACCGAAAGCATCCACAAATATTACCGGCATTATGTTGTAACAGTACCCAGCCGTGATGTGCTGA
 TTGGGGAGTTAAACAACCGATTAT

3. s4C3

probably same gene as p2D6, but different region

similarity to *S. typhimurium invA* and *Yersinia spp. lcrD*

sequence s4C3_1_U

GCGCGGACGCTAGTGTGGTGGGTGACAGCCAGACGTTACCGAACGGGATGGGGCAGATCTGTTGGCTTA
 CAAAAGACATGGCCCATAGGGCGCAAGGTTTTGGGACTGGACGTTTTCGCGGGCAGACAACGTATCTCT
 GTCTTATTAATGTGTCTGCTTCGGCATATGTATCGAACCTCGGAGCAAAGTCGTTTGGGCGCAGA
 ATTAGTACGTTTGGGTGCGTTGCTGTTATTCCTTGGGCTCGGAAAAAGAGTGCCAGCGTGAAGGAGTGG
 GATTTGGCAGACTGGCCGCCTAAT

sequence s4C3_1_R

CACTATAGGGAAAGCTTGCATGCCTGCAGGTGCGACTCTAGAGGATCTACTAGTCATATGGATTGCACTT
 GTGTATAAGAGTCAGGATTAGAGGACATGCGCCGGGAACCATATCTTTTTCCGGTGCTTCGACGCC
 ATTTGCGGAAACACAGACTTTTTGCGGCGAATGAGGATAATTGGCAATGCTAACAACGCTGAAAAGAA
 AGCGAGAGTGATAAAGGAAAGCCAGGAATTAAAGCGAGGAGCATTAAACCACAGCGGCTAATATGAG
 CGACTGAGGTTGTCTGGCAATTTG

4. p3F4

similarity to *invG* (*S. typhimurium*)

sequence p3F4_1_U

TGCAGGCCGACTCTAGAGGATCCCCGGGTACCGGTAATTTCTTTAACCTCGCATCCCCGGTGGATGAAAG
 GATATTCTGGCTGCGTAAGTAATGAATGAACCGCCAGTAGATAAAATATTGAAAGTGATAACCTGATG
 TTTTAATAACGATGCAGGATATACATATAACATGCTGGCATCAAACAGGTAAGCAAATCATATTGTGC
 TGCCAGGTTATTCAAATATCGACCGGTGCTCCAGGCGGGAATTTTCCACTAAATGTAGGTGGGATCA
 ATGGGCTAATTGGTATAGGCGGAT

Figur 6

Sheet 1 of 5

SUBSTITUTE SHEET (RULE 26)

5. p7G2

7/39

similarity to *yscC* (*Yersinia spp.*)

sequence p7G2_1_U

CCTGTGATTCCGGATGAAATAGCTTTTACGAAAGCTGTCAGACNTGCTGAAGAATACGCTGCAAATGGT
AAGCTTGTAACTTTGGGTATTGTTCCAACGCATGCTGAAACGGGTATGGATATATTCGTCGCGGTGA
GTTGATAGGAAATGACGCTTATGCAGTGGCTGAATTTGTGGAGAAACCGGATATCGATACCGCCCGTGA
CTATTTCAAATCAGGGGAAATATTACTGGCCTAGCGCGGATGTTTTATTTCGCGCAAAGCCCTTATT
AAACGAATTAACGTATCTATCACCCCAAATTCATACAGCTTGTGAA

sequence p7G2_3_0

TTACTAAACAGGGCCCCGGACCcTGTAACACCACCGcTGCCAAcACTAAAAACGATGCTTGCcGTAA
AAAAATTGAACGTTATTTACTTAATAcGCCTATTTTATTATTACATTATGCACGGACAGAGGGTGAGGATT
AAATGGATAATATTGATAATAAGTATAcTCCACAGCTATGTAAAATTTTgGGGgcTATATCgGATTcGg
TTGcTTtTAATTTAGCcTATGGcTTtCACTAGGATGTGTCTATTTTTTTtGTGGtCAAGCACAGAGAT
TTATTCcCCcAACCACC

sequence p7G2_1_I

TTTCCTTGCCGTGACAGTCCGGGATGCGAGGTTAACGAAATTACCGGCACCAAAGCTGTGGAGGTGAGC
GGTGTCcCCAGCTGCCTGACTCGTATTAGTCAATTAGCTTCAGTCTGGATAATGCGTTAATCAAACGA
AAAGACAGTGCcGGTGAAGTATATACACCGCTTAAGTATGCCACTGCGATGGATACCCAGTACCAT
TATCGCGATCAGTCCGTCTGGTTCCAGGGTCCCTAGTGTATTGCGTGAGATGAGTAACACCAGCGT
CCCGACGTCATCGACGAACAATGG

6. p9B7

similarity to *fliQ*, *invX* (*E. coli*)

sequence p9B7_1_I

CATGAGTAACCTACCCAAGTGAATCTTTACCAATATGCATCATAATCTTCTGCTGGTAAATGATTGGT
AATATCGGAAAGGTAAGTGACATAAGCACGCCATTACGTAAAAGTGCGGCCCCCTAAACTGCCACTTTTT
AATAAGGGAAGTAATAAGAAAGGCTCAATGGTCCAATAAAAGCCACAGCCAATGCAATAAGCCACTCA
TTTACCTGTTGTGCCATTCAACCATGCTCTCCAATTCGTAAACATTATCTGCCGGGTATAATTCAACAGG
ATACCGCTAAGCCATGGGTAG

sequence p9B7_3_0

ATTCAGCCCCCGGGCCATCTAACCCTATGAACAATCATCTTCTGGGTGGACAATCATTGGTACCATC
GGCCAGGCTTGTGCAATATGTATGTATCATCACGTAAAAGCGGGCCCCCTTAATCTCCCATTCTTCCTTA
AGGGCAGTTATCACGGCTGGCTCAATGGCCGGCTTAACAGCCACAG

7. s6F5

similarity to *yscU* (*Y. enterocolitica*)

sequence s6F5_1_0

GAGGCGCGTCTTCGGTTGAGGGTCGCCCTCCAGATCTTTATGCTCCTGTTTTACGTCATCTTTACTCAT
TTTAAGATCTTTTCTAATCTTATAATATTGAAAAGAATAGTCCAGTATGCCAACGACGAAATAAGAAA
CATCACCCCAACCCATAACCATTTTTTCAATGATGAAAGCACAAAGCACGCCACAGGCTAeACCACAGCC
CGGAGGGGGCCGAAAGTCTGGGATCTTGATTAATGAAAAGGCAAAGGGAAGAGATAGGATGATGCA
TGCTGGTTGGAGGCAGATTATTCATCTTCG

Figure 6

Sh et 2 of 5

SUBSTITUTE SHEET (RULE 26)

8/39

B) new sequences without similarity to entries in DNA or protein databases:**1. s4D10****sequence s4D10_1_U**

AGTTGCCGTATTTATTAAATATTCACCTCAGGTCAATATGGAGGTCTTCCCGGCTAAAAATCATTGCTT
TACTAGAGATATCACTCCCTGGGTTGCAATACAGTACGATTAGTTATCTTGATGCAGCCTGCTGATTTC
AGAATGGCAGCTGACGTACCCGCGAGACAAACATTCTGGATTATGGACGTTATCAACGCCAATATAGGG
AAGGTGGTGAAGTGGTTGATGAAATACCCCTATCCCTTGCATGTTATCGCTGACAGGACTGTTATCAGG
AGCGGGCATCCTCGATCGGCT

sequence s4D10_1_R

CAAGAGACAGATCCAACCTCGGGCCGATCGCCATAACGCCAGCAGTTTGAAAGATGAAAGCCCAGCTTAT
CCAGCCATTCCGCTACAGCGTAACGAGCAGGTTGCCAGAAATAACGATAAAGTTGCAACACCTCGGGAT
CAGGTCCGCTCAAAAACGGGCTCTCAGGCAAAAATAGCCGATCAGGATGCCCACTCCTAATAACAGTCC
TGTCACAGATAACATCAACGGATAAGGGTATTTTCATCAACCACTTCACCACCTTCCCTTTATTGGCGTT
GGATAACGTCCATAATCCAGA

2. s4H10**sequence s4H10_1_U**

AGGGCTTTTATTGATTCCATTTTACACTGATGAATGTTCCGTTGCGCTGCCCGGATTACAGCCGGATCC
TCTAGAGTCGACCTGCAGAACCGAGCCAGGAGCAAAATTAATTTTTTTGGGCAATTGCTGAAAGATGAAG
CATCCACCACTAACGCCAGTGCTTTATTACCGCAGGTTATGTTGACCAGACAAATAGATTATATGCAGT
TAACGGTAGGCGTCGATTATCTTGTGAGAATATCAGGCGCAGCATCGCAAGCGCTTAATAAGCTGGGTA
ACATGGCATGAAGGGGCAACCC

sequence s4H10_1_R

CACTATAGGGAAAGCTTGATGCCTGCAGGTCGACTCTAGAGGATCTACTAGTCATATGGATTCCTAGG
CGGCCAGATCTGATCAAGAGACAGATCCAACCTCGGGCCGATCGCCATAACGCCAGCAGTTTGAAAGATG
AAAGCCCAGCTTATCCAGCCATTCCGCTACAGCGTAACGAGCAGGTTGCCAGAAATAACGATAAAGTTG
CAACACCTCGGGATCAGGTCCGCTCAAAAACGGGCTCTCAGGCAAAAATAGCCGATCAGGATGCCCACT
CCTAATAACAGTCCTGTCAACG

3. p4G5**sequence p4G5_1_0**

CCCCCCCCCTTCTCCTGGCTTACACAGCCCCAGACCGGCGCTGGAAAAGGCCATTCCCGCCATACAGGA
GGCCAGCAACATATTTTACGCGCCGCCAGATCGTGGCCGTAACCCACGGCTTTCCGCGAGCGATTGCGC
AATCATCGCTATCGCGCCAATCGCCAGGCTGTCCGTAACGGCGTGGCGTTGAGCGCGCTGTAGGCCTC
AATCGCATGCGTCAACGCATCGATACCGGTCATCGCCGTCACGTTTGGCGGAACGCCTTCGGTCACGGA
AGCATCAAGAAATCGCCACGTCCGGC

sequence p4G5_1_U

CGCGAACGTGCCCGCAACTGCTTGTGGACGGTGAATTGCAGTTTGACGCGGCTTTCTGTCCGGAGGTC
GCCGCGCAAAAAGCGCCTGACAGCCCGCTGCAAGGCCGCGCCAACGTGATGATTTTCCCGTCGCTGGAG
GCGGGCAATATTGGCTACAAAATCACTCAGCGTCTGGGAGGCTATCGCGCTGTTGGGCGGCTAATTCAG
GGGCTTGGCGCGCGCTTCAAGACCTCTCCGAGGCTGTAGCGTCAGGAAATTATCGAACTCGCGTTG
GTGAGAAAACCAA

Figure 6

Sheet 3 of 5

4. p7A3

9/39

sequence p7A3_1_U

CGCCCTAGCATGCCTGGCGTTGTCCGGTTATTGCTCGTCAAGCGAACAGATGCAAAAGGTGAGAGCGAC
TCTCGAATCATGGGGGTCATGTATCGGGATGGTGAATCTGTGATGACTTATTGGTACGAGAAGTGCA
GGATGTTTTGGATAAAAATGGGTTACCCGCATGCTGAAGTATCCAGCGAAGGGCCGGGGAGCGTGTAA
TTCATGATGATATACAAATGGATCAGCAATGGCGCAAGGTTCAACCATTACTTGAGATATTCCCGGGT
TATTGCACTGGCAGATTAGTCACTCTC

sequence p7A3_1_I

CCCTTCCCAGGCTCGACAGGTACACAGCCAGCCACTGGTGCAGGCAGTTACTTGCTTTCATCATGGGAA
GGAGCAATATCCTGATATATTAAAGAAAGAGCGGGATCCCCCTTTCTTTACTGCTGCTAACGTTTCTTGC
AAAATGCGTTGATGAGATTCATCCAGCACACCACTGATAACAAAAGAGCGCGCATTGGCGTAACATTG
ACAAGCCCCACTAAACCGCTCTCTATTATCGCAGAAATAATATCATCCCCCTGAGACTGATGAGAGTGA
CTATTCTGCCAGCGCAAATAACCC

5. p10E11

sequence p10E11_1

ATACCGAGTATTAAGCGGCTGTGTAACATCGTCATCCAACAACATACGCAGCGAGCCGCCACGCCGGAA
AAACCGCATCGTGTATGCGCTGTTGTAGGGTCGGGTCTTTTTCATGAGTACGTTTCTGCGTATC
ATACTGGAAATTTCCCCCACTTACTGATAAGCCCTGTCAGTTGGGTAAGGACAGAGTTAAGCTCCTGA
GACATTTTTTGGAAATGGTTATCTTTCCCCGACTCATAAAATCGGTATTCCCGCTGGGGGCAATATCCAA
AGACGCTTTGGTGCCTCGCCGTAGGGCACC

sequence p10E11_U

GCGTATGCCTGCAGTTGCCCGGTTATTGCTCGTCAAGCGAACCGATGCCAAAGGTGAGAGCGACTCTC
GAATCATGGGGGTCATGTATCGGGATGGTGAATCTGTGATGACTTATTGGTACGAGAAGTGCAAGAT
GTTTTGGTAAAAATGGGTTACCCCATGCTGAAGTATCCAGCGAAGGGCGGGGAGCGTGTTAATTCAC
GATGATATTCAAATGGGTGAGCAATGGGGCAAGGTTCAACCCCACTTGCAAGATATTCCCCCCCCCTATT
GGACTGGCAGATTAGTCACTCTCA

6. s4B9

sequence s4B9_1_O

GGGCGACCTGCCCCGCGGCGCAACTTCCCCGAAGCGTTTTCCATTTCTTGTCTTAAATGACCTGGAA
AGCTTACCTAAGCCTTGTCTTGCTATGTGACAATACTGCTTGGAGAACACCCGGACGTCCATGATTAT
GCTATACAGATCACAGCGGATGGGGGATGGTGAATCGGTTATTATACCACAAGTCGCAGCTCTGAGCTT
ATTGCTATTGAGATAGAAAAACCCCCGCTTCAACTTGGATTTTGAATAATGTAATACGCAATCACCAT
ACACTATATTGCGGGTGGCGTATAA

sequence s4B9_1_R

TTGAGCTGGGGCACCGCTAATATCTTTAACCTCGCATCCCGGTGATGAAAGGATATTCTGGCTGCGTA
AGTAATGAATGAACCGCCAGCAGATAAAATATTGACAGTGATAACCCGATGTTTTTTTAACGATGCAG
GCTATACATATAACATAGCTGGCCACCAACACAGCTGAAGTAAATCATATTGTTGCTGCCAGGCTACTT
CACACTATTGTCCGGCGGGCCAGCGGGGATTTCCCCCTAAATCTCGCTGGTTCTCAAA

7. p4F8

sequence p4F8_1_I

AGTCTACGATTTTCGCTATATCTTCTCTTAATCATGGCCGCCATTTGTGGATGCGATTTTAAATATCCG
GGCGATCTTTTCATTAATAAATAAAGATTCCCCATGACTTCACAGATAAAGGTATCGGTATTTTGAGTGA
TACGTAACAATTCTGTTCTCTTCGTGTGGGTCCATGATGCGAAGAATAATGGTGGCATCATTTTCATGAG

10/39

GATTATGAACCCGAAATCTTTCTCTTTGCGATGCGCAGGCTAACTCTTCAACTCAAAAAAATCTCTG
TAAGCCGCTCTCGTGTGGGGGCGC

8. p7B8

sequence p7B8_1_O

GCGCCCCTTTAATTGGTTGAGGCGGCTGGTATTCTTGTAAGGGTAATACTAGCGAGACCCAGGTTCCAC
CCCCGGGGACACTTTTTAGTGTGAGATTACCGCCCATCATTTTAGCCAGGCTTGACGCAATAGTCAGTC
CAATTCCTGTACCTTGCGAATTTGTGTCTGCTTGATAAAAAGCAGAAAAGATTTGAGACTGCTGCTGTT
TTTCAATCCCCCACCCTATCGCTAACCAGAAATATTAATTGTTCTCCACCAAGATTGAGCGCCAGAC
GTATCCCTCCCCCTCGGGAAT

9. p8G12

sequence p8G12_1_I

GGATAAGATCCCGGATAAGTATGTCAGGCTCGTATGCACAACAGGCATTATAAACCTCTAGACCATTTT
TAACATGCTCTACTATTTTAAATGAGGCCAGGTAATAAGGCATTATAATGCCGTTAATGATGATT
CATGATCGTCTACTAATAAGATCTTATATTCTTTCAATTGGCTGCCCTCGCGAAAATTAAGATAATATT
AAGTAATGGTGTAGGTTGTGGAGATCATACGTATTTCTGGCGTAAGTCGGTTAGTTCTCCAGCGCGA
TGATTTTCCCCATTTTACGCGAT

10. p9G4

sequence p9G4_1_O.

TTCCATATTGCTCGTCCGGGAGCGTGTTAATCTTGATGATATACCAATGGATCTGCAATGGCGCAAG
GTTCAACCATTACTTTGGAGATATCCCGGGTTATTGTAAGGAGATTAGTCACTCTCATCAGTCTCAG
GGGGTGATGTTATTTCTGGGATAATAGAGCAACGGCGTTAGCAGGGGTCGGTCAGTAGTCACGGCCAA
CTTCGGTGACATTTTGCGTATCACTGGGGTATCATAACTGAATCTCATCCCCCCTTTGGTAATCAC
AC

sequence p9G4_1_U

AATCTTTTACCTCCATAAGCTGCGTGGCATAGCGATACAGAGTATTAAGCGGGTGTTACATCGTCA
TCCAAACATACGACGAGCGCCACGCGGAAAAACCGCATCGTGTATGTGCTGTTGTAGGGTC
GGGTCTTTTTTTCATGAGTACGTGTTCTGCGCTATCACTGGAATTTCCCCCCTTACTGATAAGC
CCTGTCAGTTGGGTAAGGACAGCGTTAAGCTCCTGAGACATTTTTTGAGTTGTTATCTGCCCCCGACT
CATAAGATCGGGTATTCCGCGGTGG

11. p9B6

sequence p9B6_1

ATATCCCTAATGCTTTTCTTAAATAAATACCACGGAAGGATACTGGCCACCTAGCCAAATTTAGAAA
GCAATGAACATCCGGTTTATTCTGAAAACGATTACTCCGCGCACGTTGTTCTGGCGTTACCTGAGCC
AGCAAACGATATAATGGGGTGGTGACCCGCATACCGGTCAATTGGCATCCCATCCACACCGGAGGGAGTA
AAACTCATTAGGCCATAGGTAATATCATTAAGACGCTCTAATAAATGAGGGTGGGGGCCAACTACC
ACTCCAGTATGTATTGAGTCA

12. p6G5

sequence p6G5_2_I

CCCATGGGCGCAATTTGTTGCGCAGCGTTTACCCGACCATCGCGTTTATGAGCTGTAATTCATGGGGG
TAAAAACGGCGTGACGACCCCAACGGAAGATAAGGCCGGGCTTAAACAGGAGATTATTGCTAATGCGC
AGCGCAAAGTGTTGCTGGCGGACAGCAAGTATGGCGCGCATTCGCTCTTTAATGTGGTGCCGCTTG
AGCGCTTTAATGACGTGATTACCGACGTCAATCTGCCGCCGTCAGCGCAGGTTGAACTGAAAGGGCGCG
CTTTTGGCGTAACG

Figure 6

Sheet 5 of 5

11/39

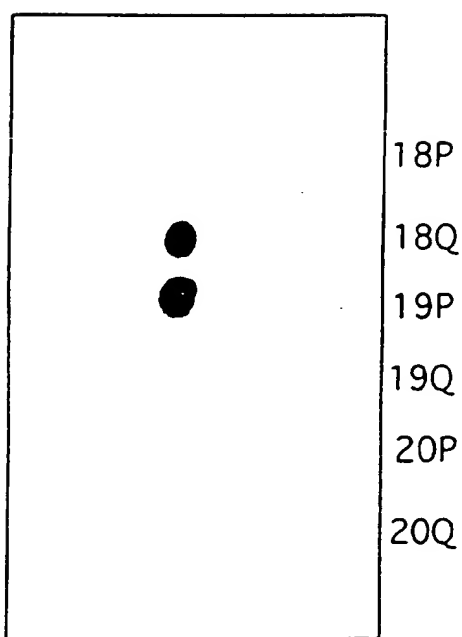


Figure 7a

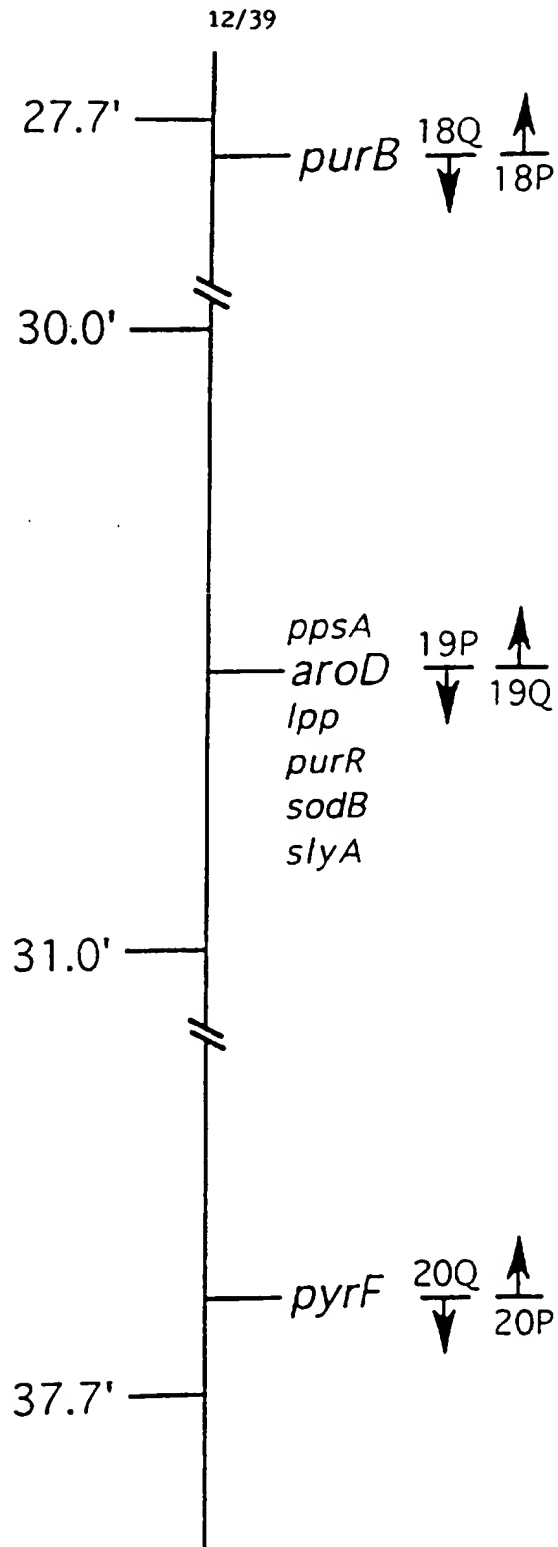
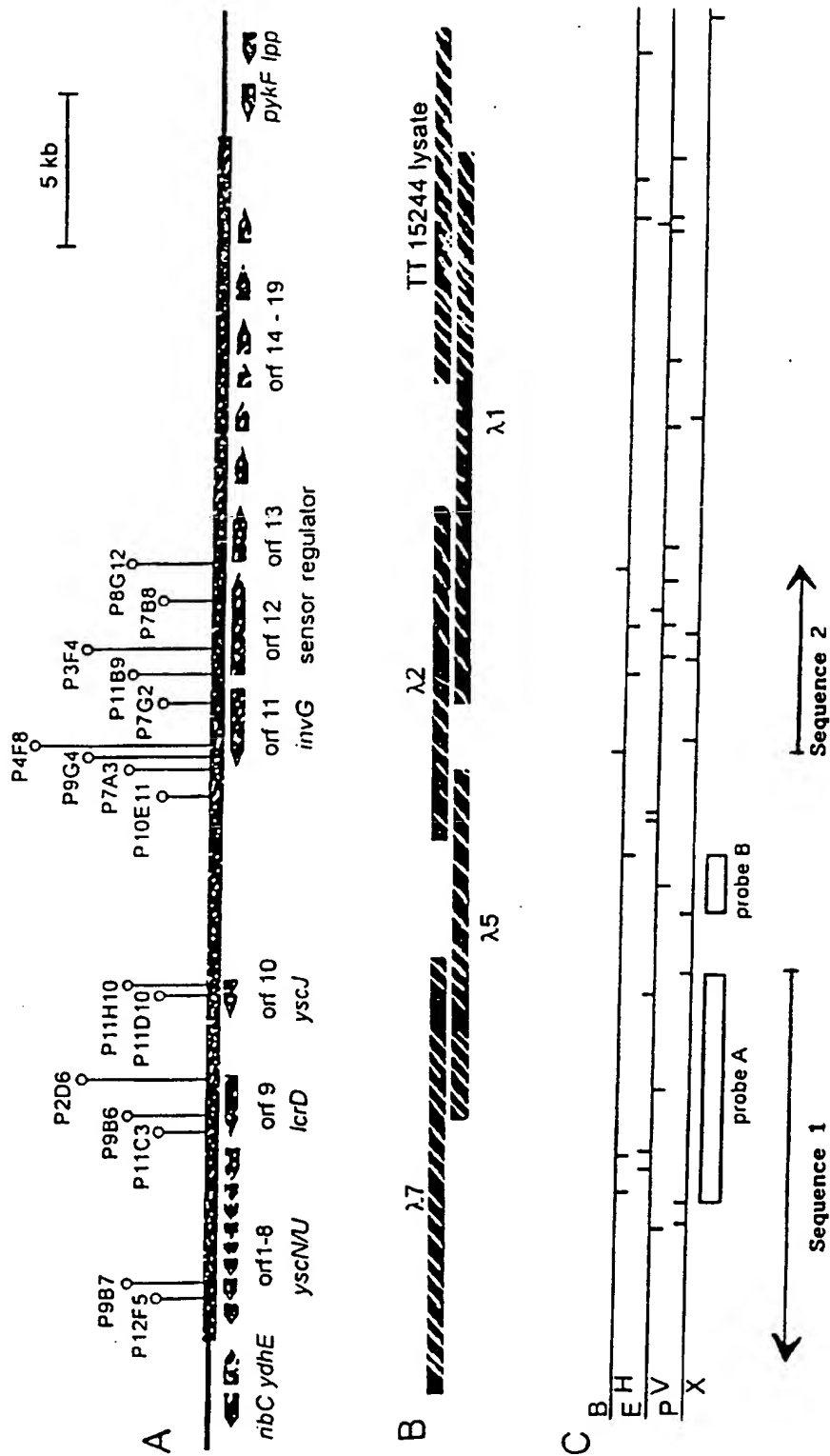


Figure 7b



Figur 8

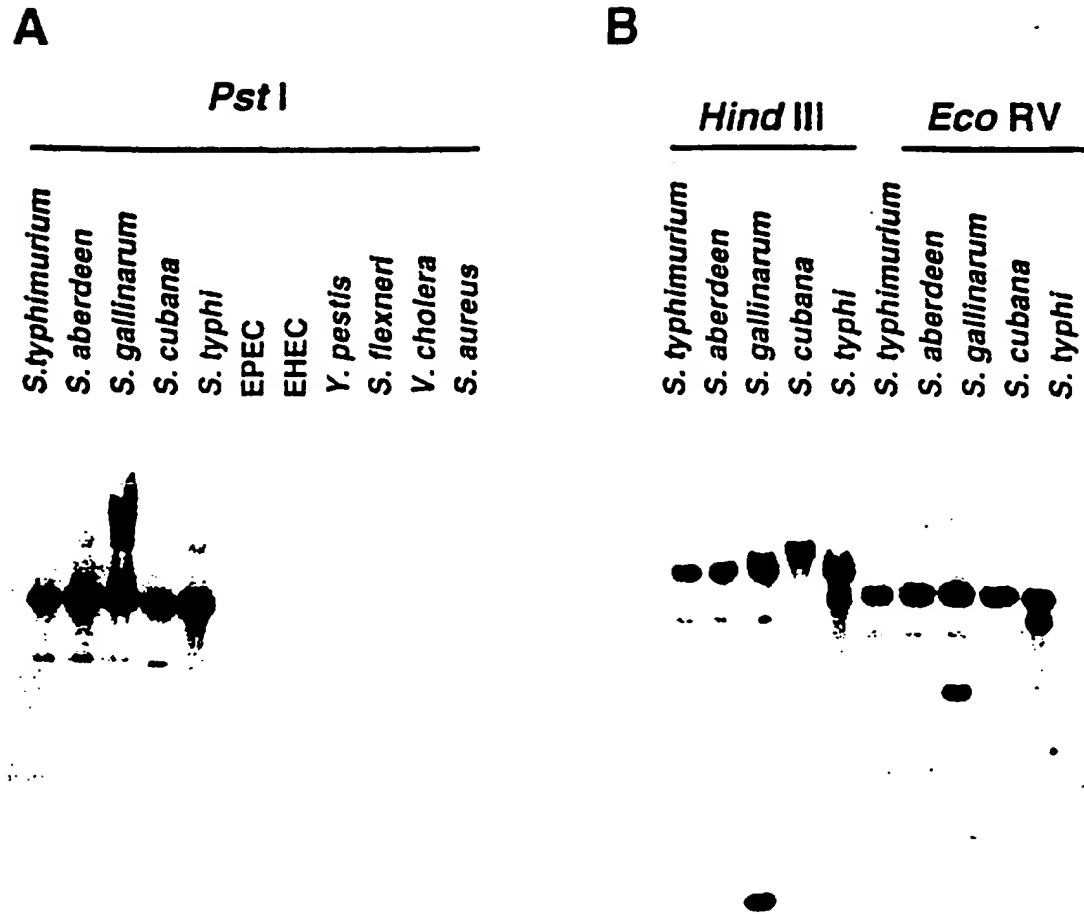


Fig. 10

16/39

DNA sequence of VGC II from centre to left hand end

CTGCAGAACCGAGCCAGGAGCAAATTAATTTTTTTGAACAATTGCTGAAAGATGAAGCATCCACCAGTAACGCCAGTGCT
1 80
GACGCTCTTGCGTCGGTCTCGTTTAATTAATAAAAACTTGTTAACGACTTCTACTTCGTAGGTGGTCATTGCGGTCACGA

L Q N R A R S K L I F L M N C * K M K H P P V T P V L -
C R T E P G A N * F F * T I A E R * S I H Q * R Q C F -
A E P S Q E Q I N F F E Q L L K D E A S T S N A S A -

TTATTACCGCAGGTTATGTTGACCAGACAAATGGATTATATGCAGTTAACGGTAGGCGTCGATTATCTTGCCAGAATATC
91 160
AATAATGGCGTCCAATACAACCTGGTCTGTTTACCTAATATACGTCAATTGCCATCCGCAGCTAATAGAACGGTCTTATAG

Y Y R R L C * P D K W I I C S * R * A S : I L P E Y H -
I T A G Y V D Q T N G L Y A V N G R R R L S C Q M I -
L L P Q V M L T R Q M D Y M Q L T V G V D Y L A R I S -

AcGGCGCAGCATGCCAAGCGCTTAATAAGCTGGATAACATGGCATGAAGGTTTCATCGTATAGTATTTCTTACTGTCCTTA
161 240
TgCGCGCTCGTACGGTTCGCGAATTATTCGACCTATTGTACCGTACTTCCAAGTAGCATATCATAAAGAATGACAGGAAT

G A A C Q A L N K L D N M A * R F I V * Y F L L S L -
T A O H A K R L I S W I T W H E G S S Y S I S Y C P Y -
R R S M P S A * * A G * H G M K V H R I V F I L T V L T -
CGTTCTTTCTTACGGCATGTGATGTGGATCTTTATCGCTCATTGCCAGAAGATGAAGCGAATCAAATGCTGGCATTACTT
241 320
GCAAGAAAGAATGCCGTACACTACCTAGAAATAGCGAGTAACGGTCTTCTACTTCGTTAGTTTACGACCGTAATGAA

start yscJ?
R S F L R H V M W I F I A H C Q K M K R : K C W H Y L -
V L S Y G M * C G S L S L I A R R * S E S N A G I T Y -
F F L T A C D V D L Y R S L P E D E A N Q M L A L L -

ATGCAGCATCATATTGATGCGAAAAAACAGGAAGAGGATGGTGTAACTTACGTGTCGAGCAGTCGGCAGTTTATTA
321 400
TACGTCGTAGTATAACTACGCTTTTTTTTGTCTTCTCCTACCACATTGGAATGCACAGCTCGTCAGCCGTCAAATAATT

start yscJ?
C S ! ! L M R K K T G R G W C N L T C R A V G S L L M -
A A S Y * C E K K Q E E D G V T L R V E Q S A V Y * -
M Q H H I D A K K N R K R M V * P Y V S S S R Q F I N -

TGCGGTTGAGGCTACTTAGACTTAACGGTTATCCGCATAGGGCAGTTTACAACGGCGGATAAGATGTTTCCGGCTAATCA
401 480
ACGCCAACTCCGATGAATCTGAATTGCCAATAGGCGTATCCCGTCAAATGTTGCCGCCTATTCTACAAAGGCCGATTAGT

R L R L L R L N G Y P H R A V Y N G G * D V S G * S -
C G * G Y L D L T V I R I G Q F T T A D K M F P A N Q -
A V E A T * T * R L S A * G S L Q R R I R C F R L I S -

GTTAGTGGTATCACCCAGGAAGAACAGGCAGAGATTAAATTTTTTAAAGAACAAGAAATGAAGGAATGCTGAGTCAG
481 560
CAATCACCATAGTGGGGTCTTCTTGTCGCTCTTAATTAATAAATTTCTTGTTCTTAACCTCCTTACGACTCAGTC

V S G I T P G R T G R R L I F * K N K E L K E C * V R -
L V V S P Q E E Q A E D * F F K R T K N * R N A E S D -
* W Y H P R K N R Q K I N F L K E Q R I E G M L S Q -

ATGGAGGGGCGTGATTAAATGGCAAAAGTGACCATTCGGCTACCGACTTATGATGAGGGAAGTAACGCTTCTCCGAGCTCA
561 640
TACCTCCCCGCACTAATTACCGTTTTCACTGGTAACGCGATGGCTGAATACTACTCCCTTCATTGCGAAGAGGCTCGAGT

W R G V I N G K S D H C A T D L * * G K * R F S E L S -
G G A * L M A K V T J A L P T Y D E G S N A S P S S -
M E G R D * W Q K * P L R Y R L M H R E V T L L R A Q -

GTTGCCGTATTATATAAATATTCACCTCAGGTCAATATGGAGGCCTTTCGGGTAAAAATTAAGATTTAATAGAGATGTC
641 720
CAACGGCATAAATATTTATAAGTGGAGTCCAGTTATACCTCCGGAAGCCCATTTTAAATTTCTAAATTTATCTCTACAG

C R I Y K I F T S G O Y G G L S G K N * R F N R D V -
V A V F I K Y S P Q V N M E A F R V K I K D L I E M S -
L P Y L * N I H L R S I W R P F G * K L K I * * R C Q -

Sequence 1

Figure 11

Sheet 1 of 19

17/39

```

      AATCCCTGGGTGCAATACAGTAAGATTAGTATCTTGATGCAGCCTGCTGAATTCAGAATGGTAGCTGACGTACCCGGCA
      7TAGGGACCCAACGTTATGTCATTCTAATCATAGAAGTACGTCGGACGACTTAAGTCTTACCATCGACTGCATGGCGCT
      800
a   N F W V A I O * D * Y L D A A C * I Q N G S * R T R E -
b   I P G L O Y S K I S I L M Q P A E F R M V A D V F A A -
c   S L G C N T V R L V S * C S L L N S E W * L T Y P R -

      GACAPACATTCTGGATTATGGACGTTATCAACGCCAATAAAGGGGAAGGTGGTGAAGTGGTTGATGAAATACCCCTTATCCG
      801
      CTGTTTGTAAAGACCTAATACCTGCAATAGTTGCGGTTATTTCCCTTCACCACTTCACCAACTACTTTATGGGAATAGGC
      880
a   T M I L D Y G R Y Q R U * R E G G E V V D E I P L S V -
b   Q T F W I M D V I N A N K G K V V K W L M K Y P Y P -
c   D K H S G L W T L S T P I K G R W * S G * * N T L I R -
      In insertion P11H1
      j

      TTGATGTTATCGTTGACAGGACTGTTATTAGGAGTGGGCATCCTGATCGGCTATTTTGCCTGAGACGCCGTTTTTGAGC
      881
      AACTACATAGCAACTGCTGACAATAATCCTACCCGTAGGACTAGCCGATAAAAACGGACTCTGCGGCAAAAACCTCG
      960
a   D V I V D R T V I R S G H P D R L F L P E T P F L S -
b   L M L S L T G L L L G V G I L I G Y F C L R R R F * A -
c   * C : R * O D C Y * E W A S * S A I F A * D A V F E F -

      CGACCTGATCCCGAGGTGTTGCAACTTTATCGTTATTTCTGGCAACCTGCTCGTTACGCTGTACCGGAATGGCTGGATAA
      961
      GCTGGACTAGGGCTCCACAACGTTGAAATAGCAATAAAGACCGTTGGACGAGCAATGCCGATGGCCTTACCGACCTATT
      1040
a   R P D P E V L Q L Y R Y F W Q P A R Y A * P E W L O Y -
b   D L : P R C C N F ! V : S G N L L V T I : R N G W ! S -
c   T * S R G V A T L S L F L A T C S L R C T G M A G * -
      In insertion P11D10
      j

      GCTGGGCTTTTCATCTTCAAACGCTGGCGTTATGGCGATCGGCCCGAGTTGGATCGTCTTCTTGACAGAGCGTTAAATAG
      1041
      CGACCCGAAAGTAGAAGTTTGACGACCGCAATACCGCTAGCCGGCTCAACCTAGCAGAGAACTGTCTCGCAATTTATC
      1120
a   L G F H L Q T A G V M A I G P S W I V F L T E R * ! D -
b   W A F ! F K L L A L W R S A R V G S S S * Q S V K * -
c   A G L S S S N C W R Y G D R P E L D R L L D R A L N P -

      ACTAAGAGGAAGCTCTGTTATTCAGCCTGTTAAATGACAGGCAAAACGGCAGGTTCGTCTTGCGCCGCTATATCGS
      1121
      TGATTCTCTTCGAGACAATAAGGTGCGACAATTTACTGTCCGTTTTTGGCGTCCAAGCGAACGCGGCGCATATAGCC
      1200
a   * E E A L L F Q P V * H T G K N G R F * L R R V Y R -
b   T K R F L C Y S S L F K * Q A K T A G S E C A A Y I G -
c   L R G S S V I P A C L N D R Q K R O V A : A P R I S A -

      CATTTGCCTTTGGGCTGGGATTATTCAAACCTCAGGTGATGACTATTTTATGCTACCAGAGTATCGGCAATTGCTTCTA
      1201
      GTAAACGGAAACCCGACCCTAATAAGTTTGAGTCCACATCACTGATAAAATACGATGGTCTCATAGCCGTTAACGAAGAT
      1280
      start 1crE?
a   H L P L G W D Y S N S G V V T I L C Y Q F : G N C F Y -
b   I C L W A G I I Q T O V * * L F Y A T R V S A I A S T -
c   F A F G L G L F N L R C S D Y F H I P E Y S Q L L L -

      CAGTGGTTTAGCGAGGATGAGATCTGGCAGCTATATGGTTGGTTGGGGCAAAGAGATGGCAAATTACTTCTCCGCAAGT
      1281
      GTCACCAATCGCTCTACTCTAGACCGTCGATATACCAACCAACCCGTTTCTCTACCGTTAATGAAGGAGGCGTTCA
      1360
a   S G L A R M R S G S Y M V G W G K E M A N Y F L R K * -
b   V V * R G * D L A A I W L V G A K R W Q I T S S A S -
c   Q W F S E D E I W Q L Y G W L G Q R D G * L L P P Q V -

      GATGCAACAACTGCATTGCAGATCGGTACCGCCATTCTTAATCGGGAAGCGCATGACGATGGGGTTTTTACATCGCGTA
      1361
      CTACGTTGTTGACGTAACGTCTAGCCATGGCGGAAGAATTAGCCCTTCGCGTACTGCTACGCCCAAATGTACCGGAT
      1440
a   C N A L H C R S V P F I L I G K R M T M * V L H A L -
b   D A T N C I A D H Y R H S * S G S A * K C G F Y M R Y -
c   M O O T A L O I G T A I L N R E A H D D A G F T C A : -

```

Figure 11

Sheet 2 of 19

18/39

TTAGTATTATTACCCCTCCGACGGTATACTTTGGCCGAAGACTTCTTTACCGAGATTATCTTCATGGAGCATTGGCT
1441 ----- 1520
AATCATAATAATGGGGGAGGCGTCGCATATGAAACCGGCTTCTGAAGAGAATGGCTCTAATAGAAGTACCTCGTAAACGA
L V L L P P F O R I L W P K T S L T E I I F M E H L L -
Y Y Y P L K S V Y F G R R L L L P A L S S W S I C Y -
S I I T F S A A Y T L A E D F S Y R D Y L H G A F A -
ATGAGTTTACTTCACTTCTCTGACGGAATTAACCATAGCTACCCGCTCGAAATATTATTGAGTCACAGTGGATAAC
1521 ----- 1600
TACTCAAAATGAAGTGAAGGAGACTGCCTTTAATTGGTATTGATGGGCGAGCTTTATAATACTCAGTGTCACCTATTG
V L L H F L R K L T I S Y P L E I L L S H S G H -
E F Y F T S S D G N P A T R S K Y Y V T V D N -
M S F T S L P L T E I N H K L F A R N I I E S Q W I T -
ATTACAATTAACCTTTATTGCGCAAGAGCAACAAGCTAAGAGAGTTTACATGCTATTGTGAGCTCCGCTTACCGTAAGG
1601 ----- 1680
TAATGTTAATTGAAATAAACCGGTTCTCGTTGTTGATTCTCTCAAAGTGTACGATAACACTCGAGGCGAATGGCATTCC
Y N L L R K S N K L R E F H M L L A P L T V R -
T I N F I C A R A T S E S F T C Y C E L R L P G -
L Q L T L F A Q E Q O A K R V S H A I V S S A Y R K A -
CTGAAAATCATCCGAGACGCTATCGTTATCAGCGTGAACAGAAAGTTGAGCAGCAACAAGAACTAGCGGTCTTGGCT
1681 ----- 1760
GACTTTTTAGTAGGCTCTGCGGATAGCAATAGTCCGACTTGCTTTCAACTCGTCTGTTCTTGTATCGCAGCAAGCGCA
L K K S S E T P I V I S V N R K L S S H N N R A C V -
N N H P R L S L S A T E S A A T F T S V L A -
E V I R D A Y R Y Q R E Q K V E Q Q Q E L A C L R -
AAAAATACGCTGGAAAAATGGAAGTGAATGGCTGGAACAGCATGTAAACATTTACAAGACGATGAAATCAATTTCC
1761 ----- 1840
TTTTATGCGACCTTTTACCTTCACCTTACCGACCTTGTCGTACATTTTGTAATGTTCTGCTACTTTAGTTAAAGC
K I R W K K W K W N G W N S H N I Y K T H K I N F V -
K Y A G P N G S G H A G T A C K T F T R S K S I S -
N T L E K M E V E W L E Q H V K H L O D D E N Q F R -
TTCATTGGTCGATCAGCGACGCATCATATTAATAATAGTATAGAACAGGTTCTGTTGGCCTGGTTCGACCAACAGTCGG
1841 ----- 1920
AAGTAACAGCTAGTGGCTCGGTAGTATAATTTTATCATATCTTGTCCAAGACAACCGGACCAAGCTGGTTGTCAGCC
H W S T Q R I L L K I V N R F C W P G S T H S R -
F I G R S R S A S Y K Y R T G S V G L V R P T V G -
S L V D A A H H I K N S I E Q V L L A W F D Q S V -
TAGACAGTGTATGTGCCATCGTCTGGCAGCCAGGCCAGGCTATGGCGGAAGAGGGAGCGCTTTATTGCGTATTAT
1921 ----- 2000
ATCTGTCAATACACGGTAGCAGACCGTGGGTCCGGTGCCGATACCGCTTCTCCCTCGCGAAATAAACCGATAAGTA
T L L C A I V W H A R P R L W R K R E A F I C F I -
R Q C Y V P S S G T P G H G Y G G R G S L F A F S S -
D S V M C H R L A R Q A T A M A E E G A L I L A I H -
CCTGAAAAGAGGCATTGATGCGAGAACTTTTGGCAAGCGGTTACGTTGATTATCGAGCCTGGTTCTCTCCCGATCA
2001 ----- 2080
GGACTTTTCTCCGTAACACGCTCTTTGAAAACCGTTGCCAAATGCAACTAATAGCTCGGACCAAGAGAGGGCTAGT
L K P R H C E K L L A S G L R L S S L V S L F I R -
K K G I D A R N F W Q A V Y V D Y R A W F L S R S -
P E K E A L M R E T F G K R F T L I I E P G F S P D O -
GGCTGAACCTTTCTCAACACGATATGCCGTTGAATTTTCACTTTCTCGTCATTTCAACGCGTTACTGAAATGGTTACGTA
2081 ----- 2160
CCGACTTGAAGGAGTGTGCTATACGGCAACTTAAAGTGAAAGAGCAGTAAAGTTGCCAATGACTTTACCAATGCAT
L N F P Q H D M P L N F H E L V I S T R I N G Y V -
G T F L N T I C R I F T F S S F O R V T E M V T -
A E L S S T R Y A V E F S L G R N R N A L L K W E R N -
ATGGTGAAAGATAAAGAGGTAGCGATGAATATTAATAATGAGATAAAATGACGCCCTACAGCATTTACCCCTGG
2161 ----- 2240
TACCACCTCTATTTTCTCATCGCTACTTATAATTTTAATTACTCTATTTTACTGCGGGGATGCTGTAATGGGGACC
M V R E V A M N I K I N E I K M T P F T A F T P G -
W R K R R I L K L M H K R P L Q H L P L A -
G E D K R G S D E Y N D K N D A P Y S I Y P W -

Figure 11

Sheet 3 of 19

19/39

CCAGGTATAGAGGAACAAGAGTTATTTCCGCTTCAATGTTAGCTCTCCAGGAGTTACAGGAAACGACGGGGCAGCGC
 2241 ----- 2320
 GGTCCATATCTCCTTGTCTCCAATAAAGCGGAAGTTACAATCGAGAGGTCTCAATGTCCTTTGCTGCCCCGTGGC
 a O V E E Q E V I S P S M L A L O E L O E T T G A A L
 b R L R N K K L F R L Q C L S R S Y R K R R G Q R
 c P G A G T R G Y F A F N V S S P G V T G N D G G S A

 TCTATGAGCAGATGGAAGAAATAGGAATGGCGCTGAGTGGTAACTGCGCGAAAATTATAAATTCAGTATGCTGAGAAA
 2321 ----- 2400
 AGATACTCTGCTACCTTCTTTATCCTTACCGCGACTCACCATTGACGCGCTTTAATATTTAAGTGACTACGACTCTTT
 a Y E T H E E I G M A L S G K L R E N Y K F T D A E K
 b S M R R W K K E W R V V N C A K I I N S L M L R N
 c L D G R N R N G A E W T A R K L I H C E T

 CTGGAGCGCAGACAGCAGGCTTTGCTGCGTTTGATAAAACAAATACAGGAGGATAATGGGGCAACGTTGCGTCCGCTTAC
 2401 ----- 2480
 GACCTCGCGCTGTGCTCGGAAACGACGCAAACTATTTTGTTATGTCCTCTATTACCCGTTGCAACGCGAGGCGAATG
 a L E R A Q A L L R L I K Q I Q E D N G A T L R P L T
 b W S L S R L C C V N K Y R R I M G O R C V R L P
 c G A Q T A G F A A F D K T H T G G W G H V A S A Y

 CGAAGAGATAGTGATCCTGATTACAGAATGCGTATCAAATATCGCTCTTGCAATGGCGCTTACTGCGCGGGGTTGT
 2481 ----- 2560
 GCTTCCTATCACTAGGACTAAATGTCTTACGATAGTTAATAGCGAGAAGTTACCGGAATGACGGCGCCCAACA
 a E E N S E F D L O N A Y Q I I A L A M A L T A G S L S
 b K R V L I Y R M R I K L S L L O W R L L P A G C
 c R R E S F T E C V S N Y R S C N G A Y C A R V V
 CAAAAAGAAAAACCGGATTTGCAATCGCAACTGGATACGTTACAGCGGAGGAGGATGGGAACCTGCCGTTTTAGTT
 2561 ----- 2640
 GTTTTTCTTTTTGCGCTAAACGTTAGCGTTGACCTATGCAATGTCGCTCTCCCTACCCTTGAACGGCAAAATCAA
 a K K K R R D L O S Q L D T L Q R R R D G N L P F L V
 b Q K R A N A I C N R N W I R Y S G G G M G T C R F F
 c K A E E T R F A I A T G Y V T A E E G W E L A V F S L

 TACTGGAATTTGGCAAGTGGATACCGTACGCTGTCTCTGGAAGCGTTTTATGCAACAGGCGATAGACACGATGAAA
 2641 ----- 2720
 ATGACCTTGACCGCTTACCTATGGCATGCGACAGGAGAGACTTCGCAAAATACGTTGTCCGCTATCTGTTGCTACTTT
 a Y W N L A K W I P Y A V L S E A F Y A T G D R Q R N
 b T G T V S G Y R T L S S L K R F M O Q A I U N D E
 c L E L S E V D T V R C P L S V L C N R E T T M H
 TGCCCTTATCGGAGTGGTTACAGCGCTGGCAGACTGGCCGGATCGCTGTGAACGGGTCGGTATTTTGTCTAAGAGCAGTA
 2721 ----- 2800
 ACGGGATACCGTACCAAGTCTGCGCACCGTCTGACCGGCTAGCGACACTTGGCCAGGCATAAAACGATTCTCGTCAT
 a A L E A V V Q T R G R L A G S L T G F F A K S S S
 b P L S Q W F R R V A D W P D R C E R V R L L R A V
 c C P E S S G S D A W Q T G R I A V N G S V F C E Q

 GCCTTTGAACTTAGCATATGCATCGAACCTCGGAGCAAAGTCGTTTGGCCGCGAGCATTAGTACGTTTGGTCTGTTTGTCT
 2801 ----- 2880
 CGGAACTTGATCGTATACGTAGCTTGGGAGCCTCGTTTCAGCAACCGGCGTCGTAATCATGCAACGCGAGCAACGA
 a L T M H M H R T L G A K S F G R S I S T F A S F A
 b A F E L S I C I E P S E Q S R L A A A L V R L R R L L
 c P L N E A Y A S N P R S K V V W P Q H Y V C V V C C
 GTTATTCCTTGGCCTTGAAAAAGAGTGCCAGCGTGAGGAGTGGATTGGCAGTTGCCGCTAATACATTACTGCCGCTAC
 2881 ----- 2960
 CAATAAGGAACCGGAACCTTTCTCACGGTCGCACTCTCACCTAACGGTCAACGGCGGATTATGTAATGACGGCGATG
 a V I P W P K R V P A G V D L P V A A Y I T A A T
 b L F L G L E K E C O R E E W I C O L P P N T L L P L L
 c Y S L A L K K S A S V R S G F A S C K L I N Y C R Y
 TACTCGATATTATTTGTGAGCGCTGGCTTTTCAGTGATTGGTTGCTTGATAGACTTACCGCTATAGTTCTTCATCGAAG
 2961 ----- 3040
 ATGAGCTATAATAACACTCGCGACCGAAAAGTCACTAACCAACGAACTATCTGAATGGCGATATCAAGAAGTAGCTTC
 a T R Y S A L A F Q L V A T Y R Y S F F I E D
 b L D I I C E R W L F S D W L L D R L T A I V S S S K
 c Y S I L F V S A G F S V I G C L I D L P L F L H R R

Figure 11

Sheet 4 of 19

20/39

```

ATGTTCAATCGGTTACTCCAACAACCTTGATGCCGAGTTTATGCTGATACCCGATAACTGTTTTAACGACGAAGATCAACG
3041 ----- 3120
TACAAGTTAGCCAATGAGGTTGTTGAAGTACGCGTCAAATACGACTATGGGCTATTGACAAAATTGCTGCTTCTAGTTGC

V Q S V T P T T * C A V Y A D T R * L F * R R R S T -
M F N K L L Q Q L D A Q F M L I P D N C F H D E D Q R -
C S I G Y S N N L M R S L C * Y P I T V L T T K I N V -

TGAACAAATTCTCGAAACGCTTCGTGAAGTAAAGATAAATCAGGTTTTATTCTGATACCTGGCTTTCAATATTAGGTAA
3121 ----- 3200
ACTTGTTTAAAGAGCTTTGCGAAGCACTTCATTCTATTAGTCCAAAATAAGACTATGGACCGAAAGTTATAAATCCATT

* T N S R N A S * S K D K S G F I L I P G F Q Y L G K -
E Q I L E T L R E V K * N O V L F * Y L A F N I * V N -
N K F S K R F V K * R * I R F Y S D T W L S I F R * -

ATTGGCTTTCTGGCTCATCATGAGGCGTCAGGATGGATTGGGATCTCATTACTGAACGTAATATTAGCTTTTATTCAA
3201 ----- 3280
TAACCGAAAGACCGAGTAGTACTCCGCGCTCTACCTAACCTAGAGTAATGACTTGCATTATAAGTCGAAATAAGTT

L A F W L I M R R Q D G L G S H Y * T * Y S A F Y S -
W L S G S S * G V R M D W D L I T E R N I Q L F I Q -
I G F L A H H E A S G W I G I S L L N V I F S F L F N -

TTAGCAGGATTAGCTGAACGGCCTTTAGCAACCAATATGTTCTGGCGGCAAGGACAATGAAACTATCATAACGGTCGT
3281 ----- 3360
AATCGTCTAATCGACTTCCGGAATCGTTGGTTATACAAGACCGCGTCTCTGTTATCTTGTATGATAGTATTGCCAGG-

S R I S * T A F S N Q Y V L A A R T ! * H * H N G R -
L A G L A E R P L A T N M F W R Q G Q Y E T I ! T V V -
* Q D * L N G L * Q P ! C S G G K D N M K L S * R S Y -

ATTCTCTTATGTCAGATACTCAAGCAACCTTCTTAGACGAAGAAGTCTTTTTAAAGCGTTGGCTAACTGGAAACCCGC
3361 ----- 3440
TAAGAGAATACAGTCTATGAGTTCGTTTGGAGAATCTGCTCTTGACGAAAAATTTCGCAACCGATTGACCTTTGGGCG

I L L C Q I L K Q T F L D E E L L F K A L A N W K P A -
F S Y V R Y S S K P S * T A N C F L K R W L T G N P Q -
S L M S D T Q A N L L R R R T A F * S V G * L E T R -

AGCGTTCCAGGGTATTCTCAACGATTATTTTTGTTGCGCGATGGGCTTGCAATGAGTTGTTCTCCACCTCTTCCAGCT
3441 ----- 3520
TCGCAAGGTCCCATAAGGAGTTGCTAATAAAACAACCGCTACCCGAACGTTACTCAACAAGAGGTGGAGAAAGTCEGA

A F Q G I P Q R L F L L R D G L A M S C S P P L S S S -
R S R V F L N D Y F C C A M G L Q * V V L H L F P A -
S V F G Y S T I I F V A R W A C N E L F S T S F Q L -

CCGCCGAGCTCTGGTTACGATTACATCATCGACAAATAAAATTTCTGGAGTCGCAATGCGTTCATGGTTAGGTGAGGGA
3521 ----- 3600
GGCGGCTCGAGACCAATGCTAATGTAGTAGCTGTTATTTTAAAGXACCTCAGCGTTACGCAAGTACCAATCCACTCCCT

A E L W L R L H H R Q I K F ? G V A M R S W I G E G -
F P S S G Y D Y I I D K * N F ? E S Q C V H G * V R E -
R R A L V T I T S S T N K I S W S R N A F M V R * G S -

GTCAGGGCGCAACAGTGGCTCAGTGTATGCCGCGGTCGCGAGGATATGGTTCTGGCGACGGTGTATTAAATCGCTATTGT
3601 ----- 3680
CAGTCCCGCGTTGTACCGAGTCACATACGCGCCAGCGCTCTATACCAAGACCGCTGCCACAATAATTAGCGATAACA

start lcrD*
V R A Q Q W L S V C A G R Q D H V L A T V L L I A I V -
S G R N S G S V Y A R V G R I W F W R R C Y * S L L * -
Q G A T V A Q C H R G S A G Y G S G D G V ! N F Y C -

GATGATGCTGTTACCTTGCCGACCTGGATGGTTGATATCCTGATTACTATCAACCTTATGTTTTAGTGATCCTGCTCT
3681 ----- 3760
CTACTACGACAATGGGAACGCGCTGGACCTACCAACTATAGGACTAATGATAGTTGGAATACAAAAGTCACTAGGACGAGA

M M L L P L P T W H V D I L I T I N L M F S V I L L L -
* C C Y P C R P G W L I S * L L S T L C F Q * S C S -
C D A V T L A D L D G * Y P D Y Y Q P Y V F S D F A L -

```

Figure 11

Sheet 5 of 19

21/39

TAATTGCTATTATCTTAGTGACCTCTCGATTTCGGTATTTCCGTCTTTATTACTTATTACTACATTATATCGTTT
3761 ----- 3840
ATTAACGATAATAGAACTACTGGGAGAGCTAAATAGCCATAAAGGCAGAAATAATGAATAATGATGTAATATAGCAAC
a I A I Y L S D P I L D L S V F P S L L L I T T L Y P L -
b L L L F I L V T L S I Y R Y F R L Y Y L L L H Y I V C -
c N C Y L S * * P S R F I G I S V F I T Y Y Y I I S F / -
TCACTCACAATCAGCACATCAGGCTGGTACTGTTACAACATAATGCCGGTAATATTGTGGATGCTTTCGGTAAGTTGT
3841 ----- 3920
AGTGAGTGTAGTCGTAGTGCAGCATGACAATGTTGTATTACGGCCATTATAACACCTACGAAAGCCATTCAACA
a S L T I S T S R L V L L Q H N A G N I V D A F G K F V -
b H S Q S A H H G W Y C Y N I M P V I L W M L S V S L S -
c T H N Q H I T A G T V T T * C R * Y C G C F R * V C -
CGTAGGAGGAATCTCACCCTGGGTTGGTTCGTATTTACCATCATTACTATCGTGCAATTTATTGTCATTACAAAGGTA
3921 ----- 4000
GCATCCTCCTTTAGAGTGGCAACCAACAGCATAAATGGTAGTAATGATAGCAGCTTAAATAACAGTAATGTTTCCAT
a V G G H L T V G L V V F T I I T I V Q F I V I T K G I -
b * E E I S P L G W S Y L P S L L S C N L L S L Q Y V -
c R R R K S H R W V G F I Y H H Y Y R A I Y C H Y R Y -
TCGAGAGGGTGGCGAAGTTAGCGCAGCTTCTCGCTGATGGGATGCCAGGCAACAAATGAGTATCGATGGCGATTG
4001 ----- 4080
AGCTCTCCACCGCTTCAATCGGTGCAAGAGCGAACTACCTACGGTCCGTTTGTCTACTCATAGCTACCGCTAAAT
a E R V A E V S A R F E L D G M P G K Q M S I D G D L -
b S R G N A K L A H V S A L M G C Q A H K * V S H A I C -
c R E G G S S * P T F L A * W D A R Q T N E Y R W R F A -
In insertion P2D6
CGTGCCGGAGTTATCGATGCAGACCATGCCCGTACATTAAGACAGCATGTCCAGCAGGAAAGCCGCTTCTCGGTGGAT
4081 ----- 4160
GCACGGCCTCAATAGCTACGTCTGGTACGGGCATGTAATTCTGTCTGACAGGTCTGCTTTCGGCGAAAGAGCCACGCTA
a R A G V I D A D H A R T L R Q H V Q Q E S R F L G A H -
b V P E L S M Q T H P V H * D S H S S R K A A F S V R W -
c C R S I R C R P C P Y I K T A C P A G K P L S R C D -
GGACGGTGGATGAATTTGTTAAAGCGGATACGATTGCCGGTATTATTGTTGTTCTGGTGAACATTATCGGGGGTATC
4161 ----- 4240
CCTGCCACGCTACTTTAAACAATTTCCGCTATGCTAACGGCCATAATAACAACAAGACCCTTGTATAGCCGCTATAG
a D G A M K F V K G D T I A G I I V V L V N I I G S I -
b T V F * N L L K A I R L P V L L L F W * T L S A V S -
c G R C C E I C * R P Y C R Y Y C C S G E H Y P A Y H -
TTATCGCTATGTAATATGATATGTCGATGAGTGAGGCTGTTACACTTATAGCGTACTGTCAATCGGAGATGTTT
4241 ----- 4320
AATAGCGATAGCATGTTATCTATACAGCTACTCACTCCGACAAGTGTAATATCGCATGACAGTTAGCCTCTACCAAT
a I A I L Q Y D M S M S E A V H T Y S V I S I G D G L -
b L S L S N M I C R * V R L F T L I A Y C Q S E M V Y -
c Y R Y A T I * Y V D E * G C S H L * R T V N R F W F H -
TGTGGGCAATTCATCGCTGCTGATTTCCCTTAGCGCGGAATTATTGTACCCCGTGTCCCGGTGAGAAAGCCAGAA
4321 ----- 4400
ACACCCGTTTAAAGGTAGCGACGACTAAAGGGAATCGCGCCTTAATAACAGTGGGCACAGGGCCCACTCTTTCGGTCT
a C G Q I E S L L I S L S A G I I V T R V P G E K R Q N -
b V G K F H R C * F P L A R L L S P V S R V R H A R -
c W A N S I A A D F P * R G N Y C H P C P G * E T F E -
CCTGGCGACAGAGTTGAGTTCTCAAATTGCCAGACAACCTCAGTCGCTCATATTAACCGCTGTGGTTTAAATGCTCCTCG
4401 ----- 4480
GGACCGCTGCTCAACTCAAGAGTTTAACGGTCTGTTGGAGTCAGCGAGTATAATTGGCGACACCAAAATTACGAGGAGC
a L A T E L S S Q I A R Q P Q S L I L T A V L M L L A -
b W R Q S * V L K L P D N L S R S Y * F L W F * T S S -
c P G D R V E F S H * Q T T S V A H I N F C G F N A P F -

Figure 11

Sheet 6 of 19

22/39

CTTTAAATTCCTGGCTTTCTTTATCACTCTCGCTTTCTTTTTCAGCGTTGTTAGCATTGCCAATTATCCTCATTCCGCCG
 4481 ----- 4560
 GAAATTAAGGACCGAAAGGAAATAGTGAGAGCGAAAGAAAGTCGCAACAATCGTAACGGTTAATAGGAGTAAGCGCGG

 L I P G F P F I T L A F F S A L L A L F : : L : R R -
 L * F L A F L L S L S L S F O R C * H C Q I S S F A A -
 F N S W L S F Y H S R F L F S V V S I A N Y P H S P Q -

 In insertion P11C3

 AAAAAGTCTGTGGTTTCCGCAATGGCGTCGAAGCACCGGAAAAAGATAGTATGGTTCCCGGGCGATGCTCTAATCTT
 4561 ----- 4640
 TTTTTCAGACACCAAGCGCTTACCAGCCTTCGTGGCCTTTTCTATCATACCAAGGGCCGCTACAGGAGATTAGAA

 K K S V V S A N G V E A P E K D S M V P G A C P L : L -
 K S L W F P Q M A S K H R K K I V W F P A H V L * S Y -
 K V C G F R K W R R S T G K R * Y G S R R H S S N L -

 ACGTCTTAGCCCGACGTTACATTCTGCCGACCTGATTCTGTATATTGACGCCATGAGATGGTTTTATTGAGGATACCG
 4641 ----- 4720
 TGCAGAATCGGGCTGCAATGTAAGACGGCTGGACTAAGCACTATAACTGCGGTACTCTACCAAAAATAAACTCCTATGGC

 R L S P T L H S A D L I R D I D A M R W F L F E D T G -
 V L A R R Y I L P T * F V I L T P * D G F I L R I P -
 T S * P D V T F C R P D S * Y * R H E M V F ! * S Y P -

 GCGTCCCTCTCCTGAGGTGAATATTGAGGTTTTGCTGAACCCACCGAAAAATTGACGGTACTGCTATATCAGGAACCG
 4721 ----- 4800
 CGCAGGGAGAGGGACTCCACTTATAACTCCAAACGGACTTGGGTGGCTTTTAACTGCCATGACGATATAGTCTTGGG

 V F L P E V N I E V L P E P T E K L T V L L Y Q E F -
 A S L S L R * I L R F C L N P P K N * R Y C * I R N P -
 R P S P * G E Y * G F A * T H R K I D G T A I S G T R -

 GTATTAGTTTATCTATTCCCGCTCAGGCGGATTATTTATTGATAGGCGCGGACGCTAGTGTGGTGGGTGACAGCCAGAC
 4801 ----- 4880
 CATAAATCAAATAGATAAGGGCGAGTCCGCCTAATAAATAACTATCCGCGCTGCGATCACACCACCCACTGTGGGTCTG

 V F S L S I P A Q A D Y L L I G A D A S V V G D S Q T -
 Y L V Y L F P L R R I I Y * * A R T L V N W V T A R R -
 I * F I Y S R S G G L F I D R R G R * C G G * Q P D -

 GTTACCGAACGGGATGGGCGAGATCTGTTGGCTTACAAAGACATGGCCCATAGGCGCAAGSTTTGGACTGGACGTTT
 4881 ----- 4960
 CAATGGCTTGGCCTACCCGCTAGACAACCGAATGTTTTCTGTACCGGGTATTCCGCGTTCCAAACCTGACCTGCAAA

 L P N G M G Q I C W L T K D M A H K A Q G F G L D V F -
 Y R T I G W G R S V G L Q K T W P I R R K V L D W T F -
 V T E R D G A D L L A Y K R H G P * G A R F W T C R F -

 TCGCGGGCAGCCAACGTATCTCTGCCTTATTAATAATGTGTCTCTCGGATATGGGAGAGTTTATTGGTGTTCAGGAA
 4961 ----- 5040
 AGCGCCCGTGGTTGCATAGAGACGGAATAATTTACACAGGACGAAGCGGTATACCTCTCAAATAACCCACAAGTCTT

 A G S Q R I S A L L K C V L L R H M G E F I G V Q E -
 S R A A N V S L P Y * N V S C F G I W E S L L V F R K -
 R G Q P T Y L C L I K M C P A S A Y G R V Y W C S G N -

 ACGCGTTATCTAATGAATGCGATGGAAAAAACTACTCTGAGCTGGTGAAAGAGCTTCAGCGCCAGTTACCCATTAATA
 5041 ----- 5120
 TCGCAATAGATTACTTACGCTACCTTTTTTGTATGAGACTCGACCACTTCTCGAAGTCGCGGTCAATGGGTAAATTAT

 T R Y L H N A M E K N Y S E L V K E L Q R Q L P I N * -
 R V I * * M R W K K T T L S W * K S F S A S Y F L I P -
 A L S N E C D G K K L L * A G E R A S A P V T H * * -

 AATCGCTGAAACTTTGCAACGGCTTGTATCAGAGCGGGTTTCTATTAGAGATTTACGTCTTATTTTCGGCACCTTAATTG
 5121 ----- 5200
 TTAGCGACTTTGAAACGTTGCCGAACATAGTCTCGCCCAAGATAATCTCTAAATGCAGAATAAAAGCGGTGAATTAAC

 I A E T L Q R L V S E R V S I R D L R L : F G T L : D -
 S L K L C N G L Y O S G F L L E I Y V L F S A P * L -
 N R * N F A T A C I R A G G F Y * M F T S : F K H L N * -

Figure 11

Sheet 7 of 19

23/39

ACTGGGCGCCACGTGAAAAAGATGTCCTGATGTTGACAGAATATGTCGGTATCGCGCTTCGTCGTCATATTCTGCGTCGT
 5201 ----- 5280
 TGACCCGCGGTGCACCTTTTCTACAGGACTACAACCTGCTTTATACAGGCATAGCGCGAAGCAGCAGTATAAGACCGCAGCA

 W A P R E K D V L M I T E Y V R I A L R R H I L R R -
 T G R H V K K M S * C * Q N M S V S R F V V I F C V V -
 L G A T * K R C P D V D R I C P Y R A S S S S Y S A S S -

 CTTAATCCGGAAGGAAAACCGCTGCCGATTTTGGCGGATCGGCGAAGGTATTGAAAACCTCGTGGCTGAATCCATTGCGCA
 5281 ----- 5360
 GAATTAGGCCTTCCTTTTGGCGACGGCTAAAACGCCTAGCGCTTCCATAACTTTTGGAGCAGCAGTATAGTAAAGCGGT

 L N P E G K P L P I L R I G E G I E N L V R E S I R Q -
 L I R K E N R C R F C G S A K V L K T S C V N P F A R -
 * S G R K T A A D F A D R R R Y * K P R A * I H S P -

 GACGGCAATGGGACCTATACTGCGCTGTCGTCTCGTCATAAGACGCAGATCCTGCAACTTATCGAGCAGGCGCTGAAGC
 5361 ----- 5440
 CTGCGTTACCCCTGGATATGACGCGACAGCAGAGCAGTATTCTGCGTCTAGGACGTTGAATAGCTCGTCCGCGACTTCG

 T A M G T Y T A L S S R H K T Q I L Q L I E Q A L K Q -
 R Q W G P I L R C R L V I P R R S C N L S S R R * S -
 D G N G D L Y C A V V S * D A D P A T Y R A G A E A -

 AGTCAGCCAAATTATTGTCACCTTCTGTGACACCCGACGTTTCTTGGGAAAAATTACAGAAGCCACCTTCTGCGAC
 5441 ----- 5520
 TCAGTCGGTTAATAAGTAACAGTGAAGACAGCTGTGGGCTGCAAAGAACGCTTTTAAATGCTTCGGTGGACACAGCTG

 S A N L F I V T S V D T R R F L R K I T E A T L F E -
 S Q P N Y S L S L L S T P D V S C E K L Q K P P C S T -
 V S Q I I H C H F C R H P T F L A K N Y R S H L V R R -

 GTACCGATTTTGTGTCAGGAAATTAGGAGAGGAGAGCCTTATACAAGTGGTAGAAAAGTATTGACCTTAGCGAGAGGA
 5521 ----- 5600
 CATGGCTAAAACAGTACCGTCCTTAATCTCTCTCTCGGAATATGTTACCACCTTTTCATAACTGGAATCGCTTCTCCT

 V P I L S W Q E L G E E S L I Q V V E S I D L S E E E -
 Y R F C H G R N * E R R A L Y K W * K V L T L A K R S -
 T D F V M A G I R R G E P Y T S G R K Y * P * R R G -

 GTTGGCGGACAATGAAGAATGAATTGATGCAACGCTCTGAGGCTGAAATATCCGCCCCCGATGGTTATTGTCGATGGGGC
 5601 ----- 5680
 CAACCGCCTGTTACTTCTTACTTAACGCTTGACAGCTCCGACTTTATAGCGGGGGGCTACCAATAACAGCTACCCCG

 end lcrD* start yscN*
 L A P N E E * I D A T S E A E I S A P R W L L S M G P -
 W R T M K N E L * M O P I L R L K Y P P P D G Y C R W G -
 V G G Q * R M N * C N V * G * N I R P P M V I V D G A -

 CGAATTCAGGATGTGAGCGCAACGTTGTTAAATGCGTGGTTGCGCTGGGGTATTATGGGCGAGTTGTGCTGTATAAGCT
 5681 ----- 5760
 GCTTAAGTCCTACAGTCGCGTTGCAACAATTTACGCACCAACGGACCCCATAAATACCGCTCAACACGACATATTTCGG

 N S G C Q R N V V K C V V A W G I Y G R V V L Y N A -
 R I Q D V S A T L L N A W L P G V F M G E L C C I N F -
 E F R M S A Q R C * M R G C L G Y L W A S C A V * S L -

 TGGAGAAGAACTTGCTGAAGTCGTGGGGATTAATGGCAGCAAAGCTTTGCTATCTCTTTTACGAGTACAATCGGGCTTC
 5761 ----- 5840
 ACCTCTTCTTGAACGACTTCAGCAGCCCTAATTACCGTCGTTTCGAAACGATAGAGGAAAATGCTCATGTTAGCCCGAAG

 yscN*
 W R R T C * S R G D * W O O S E A I S F Y E Y N R A S -
 G E E L A E V V G I N G S K A L L S P F T S T I G L H -
 E K N L L K S W G L M A A K L C Y L L L R V Q S G F -

 ACTGCGGGCAGCAAGTGATGGCCTTAAGCGACGCCATCAGGTTCCCGTGGGCGAAGCGTTATTAGGGCGAGTTATTGATG
 5841 ----- 5920
 TGACGCCCGTCGTTCACTACCGGAATTGCTGCGGTAGTCCAAGGGCACCCGCTTCGCAATAATCCCGCTCAATAACTAC

 L R A A S D G L K R R H Q V P V G E A L L G R V I D G -
 C G Q Q V M A L S D A I R F F W A K R Y * G E L L M -
 T A G S K * W P * A T F S G S R G R S V I R A S Y * W -

24/39

5921 GCTTTGGTCGTCCCCTTGATGGCCGCGAACTGCCCGACGCTGCTGGAAAGACTATGATGCAATGCCTCCTCCCGCAATG
 CGAAACCAGCAGGGGAACTACCGGCGCTTGACGGGCTGCAGACGACCTTCTGATACTACGTTACGGAGGAGGGCGTTAC 6000
 a F G R P L D G R E L P D V C W K D Y D A M P P P P A M -
 b A L V V P L H A A N C P T S A G N T M M Q C L L P O W -
 c L W S S P W P R T A R R L L E R L C N A S S R N G -
 6001 GTTCGACAGCCTATCACTCAACCATTAATGACGGGGATTGCGGCTATTGATAGCGTTGCGACCTGTGGCGAAGGGCAACG
 CAAGCTGTCGGATAGTGAGTTGGTAATTACTGCCCTAAGCGCGATAACTATCGCAACGCTGGACACCGCTTCCCGTTG 6080
 a V R Q P I T Q P L M T G I R A I D S V A T C G E G Q R -
 b F D S L S L N H R G F A L L I A L R P V A K G H E -
 c S T A Y H S T I N D G D S R Y R C D L W R R A T -
 6081 AGTGGGTATTTTTCTGCTCCTGGCGTGGGAAAAGCACGCTTCTGGCGATGCTGTGTAATGCCCGACGACGACAGCA
 TCACCCATAAAAAGACGAGGACCGCACCCCTTTTCGTGCGAAGACCGCTACGACACATTACGGCGTCTGCGTCTGCTG 6160
 a V G I F S A P G V G K S T L L A M L C N A P D A D S H -
 b W V F F L L L A W G K A R F W R C C V M R Q T Q T A -
 c S G Y F F C S W R G E K H A S G D A V C A R R R Q Q -
 6161 ATGTTCTGGTGTTAATTGGTGAACGTGGACGAGAAGTCCGCGAATTCATCGATTTTACACTGTCTGAAGAGACCCGAAA
 TACAAGACCAATAAACCCTTGCACTGCTCTTCAGGCGCTTAAGTAGCTAAATGTGACAGACTTCTCTGGGTTTT 6240
 a V L V L I G E R G R E V R E F I D F T L S E E T R A -
 b M F W C L V N V D E K S A N S S I L H C L K R P E N -
 c C S G V N W T W T R S P R I H R F Y T V R D P K T -
 6241 CGTTGTGTGATTTGTGTGCGAACCTCTGACAGACCCGCTTAGAGCGCGTGAGGGCGCTGTTGTGGCCACCACGATAGC
 GCAACACAGTAACAACAGCGTTGGAGACTGTCTGGGCGGAATCTCGCGCACTCCGCGACAAACACCGGTGGTGTATCG 6320
 a R C V E V V A T S D R P A L E R V R A L F V A T T I A -
 b V V S L L S Q P L T D P P S A G R C L W P P R G -
 c L C H C C R N L Q T R L R A R E G A V C G H H D S -
 6321 AGAATTTTTTCGGGATAATGGAAGCGAGTCTGCTTGTCTGCCGACTCACTGACGCGTTATGCCAGGGCCGACGGAAT
 TCTTAAAAAGCGCTATTACCTTTCGCTCAGCAGAACGAACGGCTGAGTGACTGCGCAATACGGTCCGGCGTGCTTTA 6400
 a E F F R D N G K R V V L L A D S L T R Y A R A A R N S -
 b N F F A I H E S E S S C L P T H R V M F G P H G N -
 c R I F S R W K A S R L A C R L T D A L C Q G R T E -
 6401 CGCTCTGGCGCGGAGAGACCGCGTTCTGGAGAATATCGCCAGGCGTATTAGTGCAITGCCACGACTTTAGAACGT
 GCGAGACCGCGGCTCTCTGGCGCAAAGACCTCTTATAGCGGTCCGCATAAATCACGTAAACGGTGCTGAAAATCTTGC 6480
 a L N F R S D R G F W R I S P G V F S A L P R L L E R -
 b R S G A G E T A V S G E Y R Q A Y L V H C H D F N V -
 c A L A F E R P R F L E N I A R R I C I A T T F R T Y -
 6481 ACGGGAATGGAGAAAAGGCAGTATTACCGCATTTTATACGGTACTGGTGAAGGCGATGATGAATGAAGCCGTTGG
 TGCCCTTACCCTCTTTCCGTCATAATGGCGTAAATATGCCATGACCACCTTCCGCTACTATACTTACTTGGCAACC 6560
 yscN
 a T G M G E A G S I T A F Y T V L V E G D D H N E A V G -
 b R E W E K K A V L P H F I R Y W W K A M I M K P I A -
 c G N G R N R Q Y Y R I L Y G T G G R R Y E S R W -
 6561 CGGATGAAGTCCGTTCACTGCTTGATGGACATATTGACTATCCCGACGCTTGACAGAGGGGGCATTATCCTGCCATT
 GCCTACTTCAGGCAAGTGACGAACCTACCTGTATAACATGATAGGGCTGCCGAACGCTCTCCCCCGTAATAGGACGGTAA 6640
 a G S F F T A W T Y C T I P T A C R E G A L S C H -
 b D E V R S L L D G H I V L S R R L A E R G H Y P A I -
 c R M K S V H C L M D I L Y Y P D G L O R G G I I L P L -

Figur 11

Sheet 9 of 19

25/39

```

GACGTGTTGGCAACGCTCAGCCGCGTTTTTCCAGTCGTTACCAGCCATGAGCATCGTCAACTGGCGGCGATATTGGGAGC
6641 CTGCACAACCGTTGCGAGTCGGCGCAAAAAGGTCAGCAATGGTCGGTACTCGTAGCAGTTGACCGCGCTATAACGCTGC
6720
a R V G N A Q P R F S S R Y Q P * A S S T G G D I A T -
b D V L A T L S R V F P V V T S H E H R Q L A A I L R R -
c T C W Q R S A A F F Q S L P A H S I V N W R R Y C D G -
GTGCCTGGCGCTTTACCAGGAGGTTGAAGCTGTTAATACGCATTGGGGAATACCAGCGAGGAGTTGATACAGATACTGACA
6721 CACGGACCGGAAATGGTCCTCCAACCTTGACAATTATGCGTAACCCCTTATGGTCGCTCCTCAACTATGTCTATGACTGT
6800
a V P G A L P G G * T V N T H W G I P A R S * Y R Y * Q -
b C L A L Y Q E V E L L I R I G E Y Q R G V D T G T D K -
c A W R F T R R L N C * Y A L G N T S E E L I Q ! L T -
AAGCCATTGATACCTATCCGGATATTTGCACATTTTTGCGACAAAGTAAGGATGAAGTATGCGGACCGGAGCTACTTAT?
6801 TTCGGTAAGTATGGATAGGCCTATAAACGTGTAAAAACGCTGTTTCATTCTCTACTTCATACGCCTGGGCTCGATGAATAT
6880
a S H * Y L S G Y L H I F A T K * G * S H R T R A T Y F -
b A I D T Y P D I C T F L R Q S K D E V C G P E L L ! -
c K P L I P I R I F A H F C D K V R M K Y A D P S Y L -
GAAAAATTACACCAAACTACTCACCGAGTGATCATGGAACCTTTGCTGGAGATAATCGCGCGGCTGAAAAGCAATTACGGC
6881 CTTTTTAATGTGGTTTATGAGTGGCTCACTAGTACCTTTGAAACGACCTCTATTAGCGCGCGGACTTTTCGTTAATGCGC
6960
end yscN* yscO*
a K I T P N T H R V I M E T L L E I I A R L E S N Y A -
b E K L H Q I I T E * S W K L C W R * S R G * K A ! T R -
c K N Y T K Y S P S D H G N F A G D N R A A E K Q L R S -
GCAAGCTTACCGTACTTGATCAGCAGCAACAGGCGGATATTACGGAACAGCAGATTTGCCAGACGCGCGCTTAGCAGTG
6961 CGTTCGAATGGCATGAAGTACTAGTCGTGTTGTCGCTAATAATGCTTGTGCTCTAAACGGTCTGCGCGGCAATCGTCAC
7040
a A S L P Y L I S S N R K L L R N S R F A R R A L * Q C -
b Q A Y R T * S A A T G D Y Y G T A D L P D A R F S S V -
c K L T V L D Q Q Q Q A I I T E Q Q I C Q T R A L A V -
TCTACCAGACTGAAAGAATTAATGGGCTGGCAAGGTACGTTATCTTGTCTTTATTGTTGGATAAGAAACACAAATGGC
7041 AGATGGTCTGACTTTCTTAATTACCCGACCGTTCCATGCAATAGAACAGTAAATAACAACCTATTCTTTGTTGTTACCC
7120
a L P D * K N * W A G K V R Y L V I Y C W I K N N K W F -
b Y Q T E R I N G L A R Y V I L S F I V G * E T T N A -
c S T R L K E L M G W Q G T V L S C H L L L D K K Q Q M A -
CGGGTTATTCACTCAGGCGCAGAGCTTTTTGACGCAACGGCAAGCAGTTAGAGAATCAGTATCAGCAGCTTGTCTCCCG
7121 GCCCAATAAGTGAGTCCGCGTCTCGAAAACTGCGTTGCGGTCGTCATCTCTTAGTCATAGTCGTCGAACAGAGGGCC
7200
a G Y S L R R R A F * R N G K O L E N Q Y Q Q L V S R -
b R V I H S G A E L F D A T A S S * R I S I S S L S P G -
c G L F T Q A Q S F L T Q R Q A V R E S V S A A C L P -
CGAAGCGAATTACAGAAGATTTTAAATGCGCTTATGAAAAAGAAAGAAAAATTACTATGGTATTAAGCGATGCGTATTA
7201 GCTTCGCTTAATGTCTTCTTAAAATTACGCGAATACTTTTTCTTTCTTTTAAATGATACCATAATTGCGTACGCATAAT
7280
end yscO* start yscP*
a R S E L Q K N F N A L M K K K E K I T M V L S D A Y Y -
b E A M Y R R I L M R L * K R K K K L L W Y * A M P I -
c K R I T E E F * C A Y E K E R K N Y Y G I K R C V L -
CCAAAGTTGAGGGAAGTCTTGGGTTGCCATGCCAGTCTTATCAGGATGATAACGAGGCGGAGGCGGAACGTATGGACTT
7281 GGTTCAACTCCCTCAGAACCAACCGGTACGGTCAGAATAGTCTCTACTATTGCTCCGCTCCGCTTGATACCTGAAA
7360
a Q S * G K S W V A M P V L S G * * R G G G G T Y G L * -
b K V E G S L G L P C Q S Y Q D D N E A E A E R M D F -
c P K L R E V L G C H A S L I R M I T R R K R N V N T L -

```

Figure 11

Sheet 10 of 19

26/39

```

GAACAACATCATGCACCAGGCATTACCCATTGGTGAGAATAATCCTCCTGCAGCATTGAATAAGAACGTGTTTTACCGCA
1361 ----- 1440
CTTGTGTGAGTACGTGGTCCGTAATGGGTAAACCACTCTATTAGGAGGACGTGTAACCTATTCTTGACCAAAAAGTGGCT

a   T T H A P G ; T H W * E * S S C S I E * E R G F H A -
b   E Q L M H Q A L P I G E N N P P A A L N K N V V F T Q -
c   N N S C T R H Y P L V R I I L L Q H * I R T W F S R N -

ACGTTATCGTGTAGTGGCGGTATCTTGACGGGTAGAGTGTGAAGTATGTGAATCAGGGGGCTAATCCAGTTAAGAA
1441 ----- 1520
TGCAATAGCACAAATCACCGCAATAGAAGCTGCCACATCTCACACTTCATACACTTAGTCCCCCGATTAGGTCAATTCTT

a   T L S C * W R L S * R C R V * S H * I R G A N P V K H -
b   R Y R V S G G Y L D G V E C E V C E S G G L I Q L P I -
c   V I V L V A V I L T V * S V K Y V N Q G G * S S * E -

TCAATGTCCCTCATCATGAAATTTACCGTTCGATGAAAGCGCTAAAGCAGTGGCTGGAGTCTCAGTTGCTGCATATGGGG
1521 ----- 1600
AGTTACAGGGAGTAGTACTTTAAATGGCAAGCTACTTTCGCGATTTCGTACCGACCTCAGAGTCAACGACGTATACCCC

a   O C P S S * N L P F D E S A K A V A G V S V A A Y G V -
b   N V P H H E I Y R S M K A L K Q W L E S Q L L H M G -
c   S M S L I M K F T V R * K R * S S G W S L S C C I W G -

TATATAATTTCCCTGGAGATATTCTATGTTAAGAATAGCGAATGAAGAGCGTCCGTGGGTGGAGATACTTCCAACGCAAG
1601 ----- 1680
ATATATTAAGGGACCTCTATAAGATACAATTCTTATCGCTTACTTCTCGCAGGCACCCACCTCTATGAAGGTTGCGTTC

end yscP* start yscQ*

a   Y N F F G D I L C * E * R M K S V R G W R Y F O R K -
b   Y I I S L E I F Y V K N S F * R A S V G G D T S N A R -
c   I * F P W R Y S M L R I A N E E R P W V E I L P T Q G -

GGGCTACCATTTGGTGAGCTGACATTGAGTATGCAACAATCCAGTACAGCAAGGGACATTATTTACCATAAATTATCAT
1681 ----- 1760
CGCGATGGTAACCACTCGACTGTAACCTCATACGTTGTTATAGGTGATGTCGTTCCCTGTAATAAATGGTATTTAATAGTA

start yscQ*

a   A L F L V S * H * V C N N I Q Y S K G H Y L P * I I I -
b   R Y H W * A D I E Y A T I S S T A R D I I Y H K L S * -
c   A T I G E L T L S M O O Y P V Q O G T L F T I N Y H -

AATGAGCTGGGTAGGGTGTGGATTGCAGAACAAATGCTGGCAGCGCTGGTGTGAAGGGCTAATTGGCACCCTAATCGATC
1761 ----- 1840
TTACTCGACCCATCCACACCTAAGCTGTTGTACGACCGTCCGACCCACACTTCCCGATTAAACGTTGGCGATTAGCTAG

a   M S W V G C G L Q N N A G S A G V K G * L A P L I D E -
b   * A G * G V D C R T M L A A L V * R A N W H R * S I -
c   N E L G R V W I A E O C W Q R W C E G L I G T A N R S -

GGCTATCGATCCTGAATTGCTATATGGAATAGCTGAATGGGGGCTGGCGCGTATTATGCAAGCCAGTGTGAACCCCTCT
1841 ----- 1920
CCGATAGCTAGGACTTAACGATATACCTTATCGACTTACCCCGACCGGGCAATAACGTTGGGTCACTACGTTGGGAG

a   L S I L N C Y M E * L N G G W R R Y C K P V M Q P S -
b   G Y R S * I A I W N S * M G A G A V I A S Q * C N P L -
c   A I D P E L L Y G I A E W G L A P L L Q A S D A T L -

GTCAGAACGAGCGCCCAACATCCTGCAGTAATCTACCACATCAGCTAGCGTTGCATATTAATGGACAGTGAAGAGCAT
1921 ----- 2000
CAGTCTTGCTCGGCGTGTAGGACGTCAATTAGATGGTGTAGTCGATCGCAACGTATAATTTACCTGTCAACTTCTCGTA

a   V R T S R Q H P A V I Y H I S * R C I L N G Q L K S M -
b   S E K A N I L Q * S T T S A S V A Y * M D S * R A -
c   Q N E P P T S C S N L P H Q L A L H I K W T V E E H -

GAGTTCATAGCATTATTTTACATGGCCAACGGGTTTTTTGCGCAATATAGTCGGAGAGCTTCTGCTGAGCGACAACA
2001 ----- 2080
CTCAAGGTATCGTAATAAAAAATGTACCGTTGCCCAAAAAACGCGTTATATCAGCCTCTCGAAAGACGACTCGCTGTTGT

a   S S I A L F L H G Q P V F C A I * S E S F L L S D N H -
b   V P * H Y F Y M A N G F F A Q Y S R R A F C * A T -
c   E F H H I I F T W P T G F L R N I V G E L S A E R Q V -

```

27/39

```

      GATTTATCCTGCCCTCCTGTGGTAGTCCCTGTATATTGAGGCTGGTGCCAGCTTACATTAAATCGAACTTGAGTCTATCG
8081 ----- 8160
      CTAATAGGACGGGGAGGACACCATCAGGGACATATAAGTCCGACCACGGTGAATGTAATTAGCTTGAACCTCAGATAGC

a      F I L P L L W * E L Y I O A G A S L H * S N L S L S -
b      N L S C P S C G S P C I F R L V P A Y I N R T * V Y R -
c      I Y P A P P V V V P V Y S G W C Q L T L I E L E S I E -

      AAATCGGCATGGGCGTTCGGATTCAATGCTTCGGCGACATCAGACTCGGTTTTTTGCTATTCAACTACCTGGGGGAATC
8161 ----- 8240
      TTAGCCGTACCCGCAAGCCTAAGTAACGAAGCCGCTGTAGTCTGAGCCAAAAACGATAAGTTGATGGACCCCTTAG

a      K S A W A F G F I A S A T S D S V F L L F N Y L G E S -
b      N R H G R S D S L L R R H Q T R F F C Y S T T W G N L -
c      I G M G V R I H C F G D I R L G F F A I Q L P G G I -

      TACGCAAGGGTGTGTGACAGAGGATAACACGATGAAATTTGACGAATTAGTCCAGGATATCGAAACGGTACTTGGCTC
8241 ----- 8320
      ATGCGTTCCACACACGACTGTCTCCTATTGTGCTACTTTAACTGCTTAATCAGGTCCTATAGCTTTGCGATGAACGCAG

a      T Q G C C * Q R I T R * N L T N * S R I S K R Y L R Q -
b      R K G V A D R G * H D E I * R I S P G Y R N A T C V -
c      Y A R V L L T E D N T M K F D E L V Q D I E T L L A S -

      AGGGAGCCCAATGTCAAAGAGTGACGGAACGCTCTTCAGTCGAACCTGAGCAGATACCACAACAGGTGCTCTTTGAGGTGG
9321 ----- 9400
      TCCCTCGGGTACAGTTTCTCACTGCCTTGCAGAGTCAGCTTGAACCTCGTCTATGGTGTGTCACAGAGAACTCCAGC

a      G A Q C Q R V T E R L Q S N L S R Y H N R C S L R S -
b      R E F N V K E * R N V F S R T * A D T T T G A L * G R -
c      G S P M S K S D G T S S V E L E Q I P Q Q V L F E V G -

      GACGTGCGAGTCTGGAATTTGGACAATTACGACAACTTAAACGGGGGACGTTTTGCCTGTAGGTGGATGTTTTGCCCA
8401 ----- 9480
      CTGACCGCTCAGACCTTTAACCTGTTAATGCTGTTGAATTTTGCCCCCTGCAAAACGGACATCCACCTACAAACCGCGT

a      D V R V W K L D N Y D E L K R G T F C L * V D V L R Q -
b      T C E S G N W T I T T T * N G G R F A C R W M F C A R -
c      R A S L E I G Q L R Q L K T G D V L P V G G C F A P -

      GAGGTGACGATAAGAGTAAATGACCGTATTATTGGGCAAGGTGAGTTGATTGCCTGTGGCAATGAATTTATGGTGGTAT
8481 ----- 8560
      CTCCACTGCTATTCTCATTACTGGCATAATAACCGGTTCCACTCACTAACGGACACCGTTACTTAATACCACGCAATA

a      R * R * E * M T V L L G K V S * L P V A M N L W C V L -
b      G D D K S K * F Y Y W A R * V D C L W Q * I Y G A Y -
c      E V T I R V N D R I : G Q G E L I A C G N E F H V R I -

      TACACGTTGGTATCTTTGCAAAAATACAGCGTAAACCTGATAAGAAAAATAATATGCGAACAATATAATAGCGTTCCAGG
8561 ----- 8640
      ATGTGCAACCATAGAAACGTTTTATGTGCAATTGGACTATTCTTTTATTATACGCTTGTATATTATCGCAAGGTCC

      end yscQ*

a      H V G I F A K I Q R K P D K K N N M R T ! * * R S R -
b      Y T L V S L Q K Y S V K L I R K I I C E Q Y N S V P G -
c      T R W Y L C K N T A * T * * E K * Y A N N I I A F Q V -

      TCGTGTGATGAGAGATACAGTATGTCTTTACCCGATTGCCTTTGCAACTGATTGGTATATTGTTTCTGCTTTCAATACT
8641 ----- 8720
      AGCACAGTACTCTCTATGTCATACAGAAATGGGCTAAGCGGAAACGTTGACTAACCATATAACAAGACGAAAGTTATGA

      start yscR*

a      S C H E P Y S M S L P D S P L Q L I G I L F L L S ! L -
b      R V M R D T V C L Y P I R L C N * L V Y C F C F O Y C -
c      V S * E I O Y V F T R F A F A T D W Y I V S A F N T -

      GCCTCTCATTATCGTCATGGGAACCTTCTTTCTTAACTGGCGGTGGTATTTTCGATTTTACGAAATGCTCTGGGTATTC
8721 ----- 8800
      CGGAGAGTAATAGCAGTACCCTTGAAGAAAGGAATTTGACCGCCACCATAAAAGCTAAAATGCTTTACGAGACCCATAAG

a      P L I I V M G T S F L K L A V V F S I I R N A L G I Q -
b      L S L S S W E L L S L N W R W Y F R F I E M L W V F -
c      A S H Y R H G N F F P * T G G G I F D F T K C S G Y S -

```

28/39

8801 AACAGGTCCGCCCCAAATATCGCACTGTATGGCCTTGGCCTTGCTACTTTCCTTATTCATTATGGGGCCGACGCTATTAGCT 8880
TTGTTTCAGGGGGTTTATAGCGTGACATACCGGAACCGGAACATGAAAGGAATAAGTAATACCCCGCTGCGATAATCGA

a Q V P P H I A L Y G L A L V L S L F I M G P T L L A -
b N K S P Q I S H C M A L R L Y F P Y S L W G R R Y * L -
c T S P F K Y R T V W P C A C T F L I H Y G A D A I S C -

8881 GTAAAGAGCGCTGGCATCCGGTTCAGGTCGCTGGCGCTCCTTTCTGGACGTCTGAGTGGGACAGTAAAGCATTAGCGCC 8960
CATTITCTCGGACCGTAGGCCAAGTCCAGCGACCGGAGGAAGACCTGCAGACTCACCTGTCATTTCGTAATCGCGG

a V K E R W H P V Q V A G A P F W T S E W D S K A L A P -
b * K S A G I R F R S L A L L S G R L S G T V K H * R L -
c K R A L A S G S G R W R S F L D V * V G Q * S I S A -

8961 TTATCGACAGTTTTTGCAAAAAAAGCTCTGAAGAGAAGGAAGCCAATTATTTTCGGAATTTGATAAAACGAACCTGGCCTG 9040
AATAGCTGTCAAAAACGTTTTTTTGAGACTTCTCTTCCTTCGGTTAATAAAAGCCTTAAACTATTTTGCTTGGACCGGAC

a Y R Q F L Q K N S E E K E A N Y F R N L I K R T W F E -
b I D S F C K K T L K R R K P I I F G I * * N E P G L -
c L S T V F A K K L * R E G S Q L F S E F D K T N L A * -

9041 AAGACATAAAAGAGGATAAAACCTGATTCTTTGCTCATATTAATTCCGGCATTACGGTGAGTCAGTTAACGCGAGGCA 9120
TTCTGTATTTTCTTTCTATTTTGGACTAAGAAACGAGTATAATTAAGGCCGTAATGCCACTCAGTCAATTGCGTCCGT

a D I K R N I K P D S L L I L I P A F T V S Q L T Q A -
b K T * K E R * H L I L C S Y * F R H L R * V S * R R H -
c R H K K K D K T * F F A H I N S G I Y G E S V N A G I -

9121 TTTGGATTGGATTACTTATTTATCTTCCCTTTCTGGCTATTGACCTGCTTATTTCAAATATACTGCTGGCTATGGGGAT 9200
AAAGCCTAACCTATGAATAAATAGAAGGGAAGACCGATAACTGGACGAATAAGTTTATATGACGACCGATACCCCTA

a F R I G L L I Y L P F L A I D L L I S N I L L A M G M -
b F G L D Y L F I R F F W L L T C L F Q I Y C W L W G -
c S D W I T Y L S S L S G Y * P A Y F K Y T A G Y G D -

9201 GATGATGGTGTCCCGATGACCATTTTATTACCGTTTAAAGCTGCTAATATTTTACTGGCAGCGGTTGGGATCTGACAC 9280
CTACTACCACAGCGGCTACTGGTAAAGTAATGGCAATTTCGACGATTATAAAATGACCGTCCGCCAACCCCTAGACTGTG

a H M V S P M T I S L P F K L L I F L L A G G W D L T L -
b * W C R R * P F H Y R L S C * Y F Y W Q A V G I * H -
c D D G V A D D H F I T V * A A N I F T G R R L G S * T -

9281 TGCGCGAATTGGTACAGAGCTTTTCATGAATGATTCTGAATTGACGCAATTTGTAACGCAACTTTTATGGATCGTCTTT 9360
ACCGCGTTAACCATGTCTCGAAAAGTACTTACTAAGACTTAACTGCGTTAAACATTGCGTTGAAAATACCTAGCAGGAAA

end yscR* start yscS*

a A Q L V C S F S * M I L N * R N L * R N F Y G S S F -
b W R N W Y R A F H E * F * I D A I C N A T F M D R P F -
c G A I G T E L F M N D S E L T Q F V T Q L L W I V L F -

9361 TTACGTCTATCCCGTAGTGTGGTGGCATCGGTAGTGTGTCATCGTAAGCCTTGTTCAGGCCTTGACTCAAATACAG 9440
AATGCAGATACGCCATCACAACCACCGTAGCCATCAACCACAGTAGCATTGCGAACAAGTCCGGAACCTGAGTTTATGTC

a L R L C R * C W W H R * L V S S * A L F R P * L K Y R -
b Y V Y A G S V G G I G S W C H R K P C S G L D S N T G -
c T S M P V V L V A S V V G V I V S L V O A L T Q I Q -

9441 GACCAAACGCTACAGTTTCATGATTAAATTATTGGCAATTGCAATAACCTTAATGGTCAGTACCCATGGCTTAGCGGTAT 9520
CTGGTTTGGATGTCAAGTACTAATTTAATAACCGTTAAGCTTATTGGAATTACAGTCGATGGGTACCGAATCGCCATA

a T K R Y S S * L N Y W O L Q * P * W S A T H G L A V S -
b P N A T V H C * I I G N C N L N G Q L P H A * R Y -
c D Q T L Q F M I K L L A I A I T L M V S Y P W L S G I -

Figure 11

Sh et 13 of 19

29/39

```

          CCTGTTGAATTATACCCGGCAGATAATGTTACGAATTGGAGAGCATGGTTGAATGGCACAAACAGGTAATGAGTGGCTTA
9521 ----- 9600
          GGACAACTTAATATGGGCCGTCTATTACAATGCTTAACCTCTCGTACCAACTTACCGTGTGTCCATTACTCACCGAAT

          end yscS*          start yscT*
a      C * I I P G R * C Y E L E S M V E W H N R * H S G L -
b      P V E L Y P A D N V T N W R A W L N G T T G K * V A Y -
c      L L N Y T R Q I M L R I G F H G * H A Q Q V N E W L I -

          TTGCATTGGCTGTGGCTTTTATTGACCAATTGAGCCTTCTTTATTACTTCCTTATTAAGGAGTGGCAGTTTATGGGGCC
9601 ----- 9680
          AACGTAACCGACACCGAAATAAGCTGGTAAGCTCGGAAAGAAATAATGAAGGGAATAATTTTACCCTCAATCCCGG

          L H W L H L L F D H * A F L Y Y F P Y * K V A V * G P -
a      C I G C G F Y S T I E P F F I T S L I K K W Q F R G R -
b      A L A V A F I R P L S L S L L L P L L K S G S L G A -
c

          GCACTTTTACGTAATGGCGTGCTTATGTCACTTACCTTTCCGATATTACCAATCATTTACCAGCAGAAGATTATGATGCA
9681 ----- 9760
          CGTGAAATGCATTACCGCACGAATACAGTGAATGAAAGGCTATAATGGTTAGTAAATGGTCGTCTTCTAATACTACGT

          H F Y V M A C L C H L P F R Y Y Q S F T S R R L * C I -
a      T F T * W R A Y V T Y L S D I T N H L P A E D Y D A -
b      A L L R N G V L M S L T F P I L P I I Y Q Q K I M M H -
c

          In insertion P987
          TATTGGTAAAGATTACAGTTGGTTAGGGTTAGTCACTGGAGAGGTGATTATTGGTTTTCAATTGGGTTTTGTGCGCGG
9761 ----- 9840
          ATAACCATTTCTAATGTCAACCAATCCCAATCAGTGACCTCTCCACTAATAACCAAAAAGTTAACCCAAAACACGCGCGC

          L V K I T V G * G * S L E R * L L V F Q L G F V R R -
a      Y W * R L Q L V R V S H W R G D Y W F F N W V L C G G -
b      I G K D Y S W L G L V T G E V I I G F S I G F C A A V -
c

          TTCCCTTTTGGGCGGTTGATATGGCGGGGTTCTGCTTGATACTTTACGTGGCGCGACAATGGGTACGATATTCATTTCT
9841 ----- 9920
          AAGGGAACCCGCAACTATACCGCCCCAAGACGAAGTATGAAATGCACCGCGCTGTACCATGCTATAAGTTAAGA

          F P F G P L : W R G F C L I L Y V A R O W V R Y S I L -
a      S L L G R * Y G G V S A * Y F T W R D N G Y D I Q F Y -
b      P F W A V D M A G F L L D T L R G A T M G T I F N S -
c

          ACAATAGAAGCTGAAACCTCACTTTTGGCTTGCTTTTCAGCCAGTCTTGTGTGTTATTTCTTTATAAGCGGCGGCAT
9921 ----- 10000
          TGTATCTTCGACTTTGGAGTGAAAAACCGAACGAAAGTCGGTCAAGAACACACAATAAAGAAATATTCCGCGCGTA

          Q * K L K F H F L A C F S A S S C V L F S L * A A A W -
a      N R S * N L T F W L A F Q P V L V C Y F L : K R R H -
b      T I E A E T S L F G L L F S O F L C V I F F I S G G M -
c

          GGAGTTTATATTAACATTTCTGTATGAGTCATATCAATATTTACCACCAGGGCGTACTTTATTATTTGACCAGCAATTTT
10001 ----- 10080
          CCTCAATATAATTTGTAAGACATACTCAGTATAGTTATAAATGGTGGTCCCGCATGAAATAATAAACTGGTCGTTAAAA

          S L Y * T F C M S H I N I Y H Q G V L Y Y L T S N F -
a      G V Y I K H S V * V I S I F T T R A Y F I I * P A I F -
b      E F I L N I L Y E S Y Q Y L P P G R T L L F D Q Q F L -
c

          TAAAATATATCCAGGCAGAGTGGAGAACGCTTTATCAATTATGTATCAGCTTCTCTCTCCTGCCATAATATGTATGGTA
10081 ----- 10160
          ATTTTATATAGGTCGCTCACCTCTTGGGAAATAGTTAATACATAGTCGAAGAGAGAAGACGGTATTATACATACCA?

          * N I S R Q S G E R F I N Y V S A S L F L P * Y V W Y -
a      K I Y P G R V E N A L S I M Y Q L L S S C H N M Y G I -
b      K Y I Q A E W R T L Y Q L C I S F S L P A I I C M V -
c

          TTAGCCGATCTGGCTTAGGTCTTTAAATCGGTCGGCACAACAATTGAATGTGTTTTCTTCTCAATGCCGCTCAAAAG
10161 ----- 10240
          AATCGGCTAGACCGAAATCCAGAAAATTTAGCCAGCGGCTGTGTTAACTTACACAAAAGAAGAGTTACGGCGAGTTTTT

          * I I W L * V F * I G R H N N * M C F S S Q C R S K V -
a      S K D D F S S F K S V G T T I E C V F L L N A A Q * K -
b      L A D L A L G L L N R S A Q Q L N V F F F S M P L N S -
c

```

Figur 11

Sheet 14 of 19

30/39

10241 TATATTGGTTCTACTGACGYCCTGATCTCATTCCCTTATGCTCTTCATCACTATTTGGTTGAAAGCGATAAATTTTATAT 10320
ATATAACCAAGATGACTGCRGGACTAGAGTAAGGGAATACGAGAAGTAGTGATAAACCAACTTTCGCTATTTAAATATA

a Y W F Y * R P D L I P L C S S S L F G * K R * I L Y -
b Y I G S T D ? L I S F ? Y A L H H Y L V E S D K F ? I -
c I L V L L T ? * S H S L M L F I T I W L K A I N F I F -

10321 TTATCTAAAAGACTGGTTTCCATCTGTATGAGCGAGAAAACAGAACAGCCTACAGAAAAGAAATTACGTGATGGCCGTAA 10400
AATAGATTTTCTGACCAAGGTAGACATACTCGCTCTTTTGTCTTGTGCGATGTCTTTTCTTTAATGCACTACCGGCATT

end yscT* start yscU*

a L S K R L V S I C M S E K T E Q P T E K K L R D G R K -
b Y L K D W F P S V * A R K Q N S L Q K R N Y V M A V R -
c I * K T G F H L Y E R E N R T A Y R K E I T * W P * -

10401 GGAAGGGCAGGTTGTCAAAGTATTGAAATAACATCATTATTTACGCTGATTGCGCTTTATTTGTATTTTCATTCTTTA 10480
CCTTCCCGTCCAACAGTTTTCATACTTTATTGTAGTAATAAAGTCGACTAACGCGAAATAACATAAAAGTAAGAAAT

a E G Q V V K S I E I T S L F Q L I A L Y L Y F H F F T -
b K G R L S K V L K * M H Y F S * L R F I C I F I S L -
c G R A G C Q K Y * N N I I I S A D C A L F V F S F L Y -

10481 CTGAAAGATGATTTTGATACTGATTGAGTCAATACTTTACATTACAATTAGTAAATAAACCATTTCTTATGCAATTA 10560
GACTTTTCTACTAAAACATATGACTAACTCAGTTATTGAAAGTGAATGTTAATCATTATTTGGTAAAAGATACGTAAT

a E K M I L I L I E S I T F T L Q L V N K P F S Y A L -
b L K R * F * Y * L S Q * L S H Y N * * I N H F L M H * -
c * K D D F D T D * V N N F H I T I S K * T I F L C I N -

10561 ACGCAATTGAGTCATGCTTTAATAGAGTCACTGACTTCTGCACTGCTGTTTCTGGGCGCTGGGGTAATAGTTGCTACTGT 10640
TGGTTAACTCAGTACGAAATTATCTCAGTGACTGAAGCGTGACGACAAAGACCCGCGACCCATTATCAACGATGACA

a T Q L S H A L I E S L T S A L L F L G A G V I V A I V -
b R N * V M L * * S H * L L H C C F W A L G * * L L L W -
c A I E S C F N R V T D F C T A V S G R W G N S C Y C -

10641 GGGTAGCGTGTCTTCTCAGGTGGGGTGGTTATTGCCAGCAAGGCCATTGGTTTTAAAAGCGAGCATATAATCGGGTAA 10720
CCCATCGCACAAAGAGTCCACCCCAACCAATAACGGTCGTTCCGGTAACCAAAATTTTCGCTCGTATATTAGGCCATT

a G S V F L Q V G V V I A S K A I G F K S E H I N F S -
b V A C F F R W G W L L P A R P L V L K A S I * : A * -
c G * R V S S G G G G Y C Q Q G H W F * K R A Y Y S G K -

10721 GTAATTTTAAGCAGATATTCTCTTTACATAGCGTAGTAGAATTATGTAATCCAGCCTAAAAGTTATCATGCTATCTCTT 10800
CATTAAAATTCGCTCTATAAGAGAAATGTATCGCATCATCTTAATACATTAGGTCGGATTTTCAATAGTAGCATAGAGAA

a N F K Q I F S L H S V V E L C K S S L K V I M L S L -
b V I L S R Y S L Y I A * * N Y V N P A * K L S C L L -
c * F * A D I L F T * R S R I M * I Q P K S Y H A : S Y -

10801 ATCTTTGCCTTTTCTTTTATTATTATGCCAGTACTTTTGGGCGCTACCGTACTGTGGGTTAGCCTGTGGCGTCTTGT 10880
TAGAAACGGAAAAAGAAATAATAACGGTCATGAAAGCCCGCGATGGCATGACACCCAATCGGACACCGCAGCAACA

a I F A F F F Y Y Y A S T E R A L P Y C G L A C G V L V -
b S L P F S F I I M P V L F G R Y R T V G * P V A C L W -
c L C L F L L L L C Q Y F S G A T V L W V S L W R A C -

10881 GGTTCCTCTTTAATAAAATGGTTATGGGTAGGGGTGATGTTTTTATATCGTCGTTGGCATACTGGACTATCTTTTT 10960
CCAAAGAAGAAATTATTTACCAATACCCATCCCCACTACCAAAAAATATAGCAGCAACCGTATGACCTGATAAGAAAAG

a V S S L I K W L W V G * M V F Y I V V G I L D Y S F O -
b F L L * * N G Y G * G * W F F I S S L A Y W T I L F -
c G F F F N K M V M G R G D G F L Y R R W H T G L F F S -

31/39

AATATTATAAGATTAGAAAAGCTATCTAAAAATGAGTAAAGATGACGTAAAAACAGGAGCATAAAGATCTGGAGGGCGACCC
 10961 ----- 11040
 TTATAATATTCTAATCTTTTCGATAGATTTTACTCATTCTACTGCATTTTGTCTCTGATTTCTAGACCTCCCGCTGG

 Y Y K I R K A I * K * V K M T * N R S I K I W R A T -
 N I I R L E K L S K N E * R * R K T G A * R S G G R P -
 I L * D * K S Y L K M S K D D V K Q E H K D L E G D P -

 In insertion P12F5

 CTCAAATGAAGACGCGCGCTCGGAAATGCAGAGTGAATACAAAGTGGGAGTTTAGCTCAATCTGTAAACAATCTGTTG
 11041 ----- 11120
 GAGTTTACTTCTGCGCGGAGCCTTTACGTCTCACTTTATGTTTACCCTCAAATCGAGTTAGACAATTTGTTAGACAAC

 L K * R R G V G N A E * N T K W E F S S I C * T I C C -
 S N E D A A S E M Q S E I Q S G S L A Q S V K Q S V A -
 Q M K T R R K K C R V K Y K V G V * L N L L N N L L -

 CGGTAGTGCCTAATCCAACGCATATTGCGGTTTGTCTTGGCTATCATCCACCGATATGCCAATACCACGCGTCTGGAA
 11121 ----- 11200
 GCCATCACCATTAGTTGCGTATAACGCCAAACAGAACGATAGTAGGGTGGCTATACGGTTATGGTGGCAGGACCTT

 G S A * S N A Y C G L S W L S S H R Y A N T T R P G K -
 V V R N P T H I A V C L G Y H P T D M P I P R V L E -
 R * C V I Q R I L R F V C L A I I P P I C O Y H A S W K -

 AAAGGCAGTGTGCTCAAGCTAACTATATTGTTAACATCGCTGAACGCAACTGCATCCCGTTGTTGAAAATGTTGAGCT
 11201 ----- 11280
 TTTCCGTCACACGAGTTGATTGATATAACAATTGTAGCGACTTGGCTTGACGTAGGGGCAACAACCTTTACAACCTCGA

 R Q * C S S * L Y C * H R * T Q L H P R C * K C * A -
 K G S D A Q A N Y I V N I A E R N C I P V V E N V E L -
 K A V M L K L T I L L T S L N A T A S P L L K M L S W -

 GGCCCGCTCATTATTTTTGAAGTGGACGCGGAGATAAAATTCCTGAAACGTTATTTGAACCCGTTGCAGCCTTGTTAC
 11281 ----- 11360
 CCGGGCGAGTAATAAAAAAATTACCTTGGCCTCTATTTAAGGACTTTGCAATAAACTTGGGCAACGTGGGAACAATG

 G P L I I F * S G T R R * N S * N V I * T R C S L V T -
 A R S L F F E V E R G D K I P E T L F E P V A A L L R -
 P A H Y F L K W N A E I K F L K R Y L N P L O P C Y -

 GTATGGTGATGAAGATAGATTATGCGCATTCTACCGAAACACCATAAATGCTTTTGGTATGCTTCTTCAGGCCACTGCGA
 11361 ----- 11440
 CATACCACTACTTCTATCTAATACGCGTAAGATGGCTTTGTGGTATTTACGAAAACCATACGAAGAAGTCCGGTGACGCT

 end *yscU*
 Y G D E D R L C A F Y R N T I N A F G M L L Q A T A * -
 M V M N I D Y A H S T E T P * M L L V C F F R F L R -
 V W * * R * I M R I L P K H H K C F W Y A S S G H C E -

 AGGTAAAGAGGGTAATAGCGTATAGAGCAGTGCTTGACGATAAAGGTGAGAGACTGAAAATAATCGCTTTTAGCCTGGCA
 11441 ----- 11520
 TCCAATTCTCCATTATCGCATATCTCGTCACGAACTGCTATTTCCACTCTCTGACTTTTATTAGCGAAAATCGGACCGT

 V K R V I A Y R A V L D D K G E R L K I I A F S L A -
 R L R G * * R I E O C L T I K V R D * K * S L L A W H -
 G * E G N S V * S S A * R * R * E T E N N R F * P G T -

 CAAGCACCAGATAGCGTATTATAAAATTAACAAGATAATGGATTGGTGGTCTGAATGGACTCGAACCACTCGACCCCC
 11521 ----- 11600
 GTTCGGTGGTCTATCGCATAATATTTAATTGTCTATTACCTAACCCACGAGACTTACCTGAGCTTGGTgAGCTGGGGG

 Q A P D S V L * N * T R * W I G A S E W T R T T R P P -
 K H Q I A Y Y K I K Q D N G L V R L N G L E P L D P H -
 S T R * R I I K L N K I M D W C V * M D S N H S T P -

 ACCATGTCAAGGTGGTGGTCTAACCACCTGAGCTATGAACGGCAACGTTGTAGGTGACAACGGGGACGAATATTAGCGTC
 11601 ----- 11680
 TGGTACAGTTCCACCACGAGATTGGTTGACTCGATACTTGCCGTGCAACATCCACTGTTGCCCTGCTTATAATCGCAG

 P C O G G A L T N * A M N G N V V G D N G D E Y * R H -
 H V K V V L * P T E L * T A T L * V T T G T N I S V -
 T M S R W C S N Q L S Y E R Q R C H * Q R G R I L A S -

32/39

```

      ACAACCGCAATGAGGCAAGAGGGAATCGCAATTTTCTTCTGAAATCACCTGATTGCGGTGGAAATATGCAACATGTCTG
11691 ----- 11760
      TGTGGCGTTACTCCGTTCTCCCTTTAGCGTTAAAGAAGGACTTTAGTGGACTAACGCCACCTTTATACGTTGTACAGC

a      N R N E A R G K S Q F S S * N H L I A V E I C N M S -
b      T T A M R Q E G N R N F L P E I T * L R W K Y A T C R -
c      Q P Q * G K R E I A I F F L K S P D C G G N M Q H V E -

      AGAAAATAGCCGCATGCGACGGCTATCGTCGTATTATCGGAGCGCGCTGCAAAATGATGGCGGACGGCTGACGTTGTAG
11761 ----- 11840
      TCTTTTATCGGCGGTACGCTGCCGATAGCAGCATAATAGCCTCGCGGACGTTTTACTACCGCTGCCGACTGCAACATC

a      R K * P P C D G Y R R I I G A R C K M M A D G * R C R -
b      E N S R H A T A I V V L S E R A A K * W R T A D V V D -
c      K I A A M R R L S S Y Y R S A L Q N D G G R L T L * -

      ATAGCGCATCCGTAGCATCATTAAACACCGCGCGGAGGTGAGGCGGATGATGAACCCATCCAGAAGCCTGCCGTTCCCA
11841 ----- 11920
      TATCGCGTAGGCATCGTAGTAATTGTGGCGGCGGCTCCAGTCCGGCTACTACTTGGGGTAGGTCTTCGGACGCCAGGGT

a      * R I R S I I N T A A E V R P M H N P I O K P A G P I -
b      S A S * S L T P P P R S G R * * T P S R S L P V P -
c      I A H F * H H * H R R R G Q A D D E P H P E A C R S H -

      TACGATCCACCACCAATCCGTTAACGCCAGGATATAACCGCTGGGTAAACCTAACACCCAGTAGGCGGTAAAGGTGATA
11921 ----- 12000
      ATGCTAGGTGGTGGTTAGGCAATTGCGGTCTTATATTGGCGACCCATTGGATTGTGGGTCATCCGCCATTCCACTAT

a      R S T T N S V N A R I * P L G K P N T O * A V K V I -
b      I D P P N P L T P G Y N R W V N L T P S R R * R * * -
c      T I H H Q I R * R Q D I T A G * T * H P V G G K G D K -

      AAAAAATGGAACGCGTATCTTTATAACCGCGCAGAATACCGCTGCCGATAACCTGTATAGAGTCGGAAATCTGGTAAAC
12001 ----- 12080
      TTTTCTACCTTGCGCATAGAAATATTGGCGCGTCTTATGGCGACGGCTATTGGACATATCTCAGCCTTTAGACCATTGT

a      K K M E R V S L * P R R I P L P I T C I E S E I W * T -
b      K R W N A Y L Y N R A E Y R C R * P V * S R K S G K P -
c      K D G T R I F I T A Q N T A A D N L Y R V G N L V N -

      CGCAGCGAGCAGCATTAATTGCGGCAAGCGCCACGACCTCAGGGTTGTCATGTAGAGCAAAGCAATATGCTTACGCAGA
12081 ----- 12160
      CGCTCGCTCGTCGTAATTAACGCCGTTGCGGTGCTGGAGTCCCAACAGTAACATCTCGTTTCGTTATACGAATGCGTCT

a      A A S S I N C G K R H D L R V V I V E Q S N M L T Q S -
b      Q R A A L I A A S A T T S G L S L * S K A I C L R R -
c      R S E Q * L R Q A P R P Q G C H C R A K Q Y A Y A E -

      GTAACGGTAAAAATAGCGGTAAACACAGCCATACAAATGCCGACGCCCTAAACCGGTACGCGCTCGGTTTGCGCATCCAGC
12161 ----- 12240
      CATTGCCATTTTATCGCCATTGGTGTGCGGTATGTTACGGCTGCGGATTGGCCATGCGCGACGCAACGCGTAGGTGC

a      N G A N S G N H S H T N A D A * T G T R C V C A S S -
b      V T V N I A V T T A I Q M P T P K P V R A A F A H P A -
c      * R * K * R * P Q P Y K C R R L N R Y A L R L R I Q R -

      GTTGAGCCCTGGCCAGACCGATAACCCACTCGAATCGTTACCGCGCAGCCAGCGACATCGGCAGTACGAACATCAGCG
12241 ----- 12320
      CAACTCGGGACCGGCTCTGGCTATTGGGTGAGCTTAGCAATGGCGGCGTGGTGGCTGTAGCCGTCATGCTTGTAGTCGC

a      V E P W P R P I T H S N R Y R R S Q R H R Q Y E H Q R -
b      L S F G P D R * P T R I V T A A A S D I G S T N I S E -
c      * A I A Q T D N P I E S L P P O P A T S A V R T S A -

      AGCTAAAGTTAAGCGCAATCTGATGACCGCGGACATCCACAATACCTAATGGCGAAACAGCAGCGCAACGACCGCAAT
12321 ----- 12400
      TCGATTTCAATTCCGCTTAGACTACTGGCCGCTGAGGTGTTATGGATTACCGCTTTGGTCGTGCGGTTGCTGGCGTTTA

a      A K V K R N L M T G D I H N T * W R N Q Q R N D R K * -
b      L K L S A I * * P A T S T I P N G E T S S A T T A N -
c      S * S * A Q S D D R R H P Q Y L M A K P A A Q R P Q I -

      AACGTCACTTCAGAGAACGCGCAATCGGCAACCCAGTTGAATCAGGCGCTTCATGACGACGCTATCGGGTTTGC
12401 ----- 12480
      TTGCAGTGAAGTTCTTGTGCGGTGCGGTAGCCGTTGGGGTCAACTAGTCCGGAAGTACTGCTGCGATAGCCCAAACG

a      R H F K E Q P A Q S A T P V E S G A S * R R Y R V C -
b      N V T S Y N S Q R N R O P Q L N Q A L H D D A I G F A -
c      T S L Q H T A S A I G N P S * I R R F M T T L S G L P -

```

Figure 11

Sheet 17 of 19

33/39

CAAAGCCTTTTTCATTACGAATATCACGCATTGAACGCGCGTGTAAATGTAAGAAAGCATGGCGATAAACATCACCCAA
 12481 ----- 12560
 GTTTCGGAAGTAATGCTTATAGTGCCTAAGTTCGCGCACAAATTACATTCTTTCGTACCGCTATTTGTAGTGGGT
 Q S L F H Y E Y H A L N A R V . C K K A W R . T S P N -
 K A I F I T N I T H . T R V F N V R K H G D K H H P I -
 K P F S L R I S R I E R A C L M . E S M A I N I T Q -
 TAGACCGCGCAGTCGCAACGCGCGAGCCGATACCGCGAGTTCGCGCATACCAAAATGGCCATAGATAAAAAATATAGTT
 12561 ----- 12640
 ATCTGGCGGCGTCAGCGTTCGCGCGTTCGCGTATGGCGGCTCAAGGCGGTATGGTTTACCGGTATCTATTTTATATCAA
 R P P Q S Q R R S R Y R R V P A Y Q N G H R . K Y S S -
 D R R S R N A A A D T A E F R H T K M A I D K N I V -
 T A A V A T P Q P I P P S S G I P K W P . I K I . F -
 CACCGGAATATTCACCGAGCGGCCAAAAATCCCATCACCATACCGGTTTGGTTTGGCCAGACCTTCGCACTGGTTTC
 12641 ----- 12720
 GTGGCCTTATAAGTGGTCGTCGGGTTTTAGGGTAGTGGTATGGCCAAACCAAAACCGGTCTGGAAGCGTGACCAAAG
 P E Y S P A G P K I P S P Y P V W F W P D L R T G F -
 H R N I H Q Q A Q K S H H T R F G F G Q T F A L V S -
 T G I F T S R P K N P I T I P G L V L A R P S H W F R -
 GCGCTACCTGAAGAAAGGTATCCTGCGCCCCACAGCAGCGCGCGGAAGATAACCCACGGCTTTATCGGCCAGCGCGGGA
 12721 ----- 12800
 CGCGATGGACTTTCTTTCCATAGGACGCGGGGTGTCGTCGCGCGCTTCTATTGGTGCCGAAATAGCGGTCGCGGCT
 A L P E R K G I L R P T A A R E D H P R L Y R P A P D -
 R Y L K E K V S C A P Q Q R A K I T H G F I G Q R R I -
 A T . K K R Y P A P H S S A R R . P T A L S A S A G -
 TCAATATTATGCATAGAGCGGATAATGTATCCGGCATTCCACAGGACGATCATCACCGACGAGACAAAGCCCGCCAG
 12801 ----- 12880
 AGTTATAATACGTATCTCGCCTATTACATAGCGCGTAAGGTGTCTGCTAGTAGTGGTCGTGCCTCTGTTTCGGGCGGTC
 Q Y I A . S G . C I R H S T G R S S P A R R Q S P P A -
 N I M H R A D N V S G I P Q D D H H Q H G D K A R Q -
 S I L C I E R I M Y P A F H R T I I T S T E T K P A S -
 CCAGAACCCTTGTCGAACCTGATGCGCGATACGCTCACGACGCGCGGAGCCATTGAGTTGCGCAATCACAGGCGTCAAGG
 12881 ----- 12960
 GGTCTTGGAACAGCTTGGACTACGCGCTATGCGAGTGCTGCGGCGCTCGGTAACCAACGCGTGTAGTGTCCGCGAGTTCC
 R T L V E P D A R Y A H D G R S H . V A Q S Q A S A -
 P E F L S N L M R D T L T T A G A I E L R N H R R Q S -
 Q N F C R T . C A I R S R R P E P L S C A I T G V K A -
 CCAGCAGTAAGCGGTGACCAAAATGGCGGGAAGCAGATAGAGGTGCCGATAGCGAGCGCAGCCATGTCCGTAGCGG
 12961 ----- 13040
 GGTGCTCATTCGGCACTGGTTTGTTTACCGCCCTTCGTCTATCTCCACGGCTATCGCTGCCGTGGTACAGGCATCGCG
 P A V S R D Q T K W R E A D R G A D S D G S H V R S A -
 Q Q . A V T K Q N G G K Q I E V P I A T A A M S V A L -
 S S . P . P N K M A G G S R . R C R . R R Q P C P . R -
 TATAGCCTCCCGCATGACGGTATGACGAATCCATTGCGGCTATACCACTTGGCGAAGGATCACCGGTATCTGAACGC
 13041 ----- 13120
 ATATCGGAGGCGGTAAGTACTGCTAGTAAAGCGGATATGGTGAACGCGTTCCTAGTGGCCATAGACTTGGC
 I A S R H D G I D E S I A V Y T T C A R ! T G I . T L -
 . P P A M T V S T N P L R S I P L A Q G S P V S E R -
 Y S L P P . R Y R R I H C G L Y H L R K D H R Y L N A -
 TAATAACTGACGCGCTTCACTGGTATACTTCTGCAGTATTACCTTTTATTTTGTGTATATGAAAGACTAAAAAGCC
 13121 ----- 13200
 ATTATTGACTGCGGAAGTGACCATATGAAGACGTGCATAAGTGGAAAATAAAACAATATACTTTCTGATTTTTCGG
 I T D A L H W Y T S A R I H L L F C C Y H K D . K A -
 . L T R F T G I L L H V F T F Y F V V I . K T K K P -
 N N . R A S L V Y F C T Y S P F I L L L Y E R L K S P -
 GCCGAAGTGGCAGCCAAAAGAAATAGCAGGGGAAATTCAGTCTATTGTAGCGGGGTATTACTATTCTCCAGTGAAAAA
 13201 ----- 13280
 CGGCTTACCGTCGGTTTTCTTTATCGTCCCCTTAAAGTCAGATAACATCGCCCCATAATGATAAAGAGGTCACTTTTT
 A E V A A . R N S H G N I S L L . R G I T I S P V K K -
 P K W Q P K E I A G E I S V Y C S G V L L F L Q . K N -
 R S G S Q K K . Q G K F O S I V A G Y Y Y F S S E K -

Figure 11

Sheet 18 of 19

34/39

```
13281 ACAGTTGTTAACGGCGCATTGCTGGCAAGCTGTTTTCCACCTGCTATTGTGCTGAACAGTTCTGCTTTTATTATTCA
----- 13360
TGTCACAATTGCCGCGTAACGACCGTTCGACAAAAAGGTGGACGATAACACGACTTGTCAAGACGAAAAATAAATAAGT
a      Q L L T A H C W Q A V F P P A I V L N S S A F I Y F R -
b      S C * R R I A G K L F F H L L L C * T V L L L F I S -
c      T V V N G A L L A S C F S T C Y C A E Q F C F Y L F Q -

13361 GGAGTTGAAGATATGTTTACGGGGATCGTACAGGGTACCGCGAAACTGGTATCGATA
----- 13417
CCTCAACTTCTATACAAATGCCCTAGCATGTCCCATGGCGCTTGACCATAGCTAT
a      S * R Y V Y G D R T G Y R E T G I D -
b      G V E D M F T G I V Q G T A K L V S I -
c      E L K I C L R G S Y R V P R N W Y R -
```

35/39

DNA sequence of VGC II cluster C

Tn insertion P9B4

GGATCCTTTTCTTTAATGCTGCTAACGTTTCTTGCAAAATGCGTTGATGAGATTATCCAGTACACCACTGATAACAAA
CCTAGGAAAAAGAAATTACGACGATTGCAAAGAAGCGTTTACGCACTACTCTAAGTAGGTCATGTGGTGACTATTGTTT

Tn insertion P7A3

AGAGCGNCGCATTGGCNWMMHTKRNNMRNNSCNNACTAAACCGTTCTCTATTATCGCAGAAATAATATCATCCCCCTG
TCTCGNCGCTAACCGNHTKKWAMYNKNYNSGNNTGATTGGCAAGAGATAATAGCGTCTTTATTATAGTAGGGGAC
AGACTGATGAGAGTGACTAATCTGCCAGTGCAATAACCCGGGAATATCTGCAAGTAATGGTTGAACCTTGCGCCATTGCT
TCTGACTACTCTCACTGATTAGACGGTCACGTTATTGGGCCCTTATAGACGTTCACTACCAACTTGGAAACGCGGTAACGA

Tn insertion P96A

GATCCATTTGTATATCATCATGAATTAACACGCTCCCCGGCCCTTCGCTGGATACTTCAGCATNSSGGTAACCCATTTTT
CTAGGTAAACATATAGTAGTACTTAATTGTGCGAGGGGCGGGAAGCGACCTATGAAGTCGTANSSCCATTGGGTAAAAA
ATCAAAACATCCTGCACCTTCTCGTACCAATAAGTCATCACAGATTACACCATCCCGATACATGACCCCCCATGATTCGAG
TAGTTTTGTAGGAAGTGAAGAGCATGGTTATTTCAGTAGTGTCTAATGTGGTAGGGCTATGTACTGGGGGGTACTAAGCTC
AGTCGCTCTCACCTTTTGCATCTGTTGCTTGACGAGCAATAACCGGACAACGCAGGCTGCCATCTTCTTTCCATTGGG
TCAGCGAGAGTGGAAAACGTAGACAAGCGAACTGCTCGTTATTGGCTGTTGACGTCGACGGTAGAAGAAAGGTAAACGC
CCCGCACATAATGAATATTGCTTTTGTCTAATAAAAACTTAACCCGCAAAGGTAAGTCATTTACCGTTTCAGGCTGACCA
GGGCGTGATTACTTATAACGAAAACAGATTATTTTGAATTGGGCGTTTCCATTAGTAAATGGCAAAGTCCGACTGGT
CTAATACTTAACAGGACACCCATTCCACCGATGAAATCAAGAATACGCCAGCCAACCACAGTACCCTGATCTGGAAC
GATTATGAATTGCTCTGTGGTAAGGTGGCTACTTTAGTTCTTATGCGGTGGTGGTGGTATGGGACTAGACCTTTG
GGGTATTTGATAATCAGCAAGTTCACAATCCTGTTTACCAAACGCGATASSCACTCCCGCAACCTGCAAAACCCCACTGG
CCCATAACTATTAGTCGTTCAAGTGTAGGACAAATGGTTTGCCTATSSGTGAGGGCGTTGGACGTTTGGGGTGACC
ATGGTAGCGGCTTTTGGATTAAATCTGCGGCCATTAACCTAATCTGGCTTTCCCGGCATCAACAAATAAATCTATCT
TACCATCGCCGAATAAACCTAATTTAGACGCCGTAATTGAGATTGAGACCGAAAGGGCGTAGTTGTTTATTGATAGA
GCCTGTTCTCTCAGAATAATTTTTCATTATAGCCAGCGAATACAAATATCGCATCCCTTCTCCCCAGTGACAGGTTA
CGGACAAGAGAGTCTTATTAATAAAGTAAATATCGGTGCTTATGTTTATAGCGTAGGGAAGAGGGGGTCACTGTCCAAT
CCTTCATTGAGCCATCTTCCGGGCTTGTAACAGTGACCTAAAAACGTATTTTCCAGGAACCTTTTGGAATTAACCAT
GGAAGTAAGTCGATGAAGGGCCGGAACATTTGCACTGGATTTTTCATATAAAGGTCCTTGAGAAACCTAATTGGTA
GAGATATGCCATTATTTACTACTGAGGCTTAAATCAAAAAAGCCTGATTACACTATGTACTTGAGTCGTATCATTGCGA
CTCTATACGGTAATAATGATGACTCCGAAATAGTTTTTTTCGGACTAATGTGATACATGAACTCAGCATAGTAACGCT
AACAATGACCTACACAGGAATATCGCCCAATAAGGGATTTGTTTTGCGAGTGGATTGTTTACCTTGTTTAAACCC
TTGTTTACTGGATGTTGCTTATAGCGGGTATTTCCCTAAAAACAAACGCTCACCTAAACAAATGGAACAAATTTGGG
TCCAGCAATNAGACTTTGCCCGGCCAATAATGTGGCTTGCGAANCRAATTCAGAAATTTGCACTTCGGGCAGCGGGTCT
AGGGTCGTTANTCTGAAACGGGCGGTTATTACACCGAACGCTTNGYTAAGTCTTAAACGTGAAGCCCGTCGCCGAGA
GTNTYGCYTTKGNSTATCACTTTGTTGTCCATCCTGAANTATTAAGATTAAGCATTATTTTTTTCGCGCCATTGTCAATTT
CANARCGRAAHCHNSATAGTGAAACACAGGTAGGACTTNATAATTCTAATTCGTAATAAAAAACGCACGGTAACAGTAAA
AACAAGCGAGGTGTAACCGGNNAACAAAGAACCCGTAGTGATGGATTCAAGTTTAGCCACTTTTTCTCCCTGCAGTTTGG
TTGTTTCGCTCCACATTGCGCHNTTGTCTTGGGCATCACTACCTAAGTTCAAATCGGTGAAAAAGAGGGACGTCAAACC

Sequence 2

Figure 12

Sheet 1 of 5

36/39

1361 TATAGAAAGTAATATTTTTATCCAGCACAGCCTGGATATTATTTAAAGTCACCACAGATGGCTGGGAAAGTACATAAGCC
 ATATCTTTTATTATAAAAAATAGGTCGTGTGGACCTATAATAAATTCAGTGGTGTCTACCGACCCTTCATGTATTCCG 1440
 1441 TGAGAGCTTTTTTCCAGGGCATTTCAGACGCACCATAAAGTTTGAGGTATCGCTGATTACCGTTGANNNAACCACTAGCACC
 ACTCTCGAAAAAAGGTCOCGTAAAGTCTGCGTGGTATTTCAAACCTCCATAGCGACTAATGGCAACTNNTTGGTGATCGTGG 1520
 1521 ACCGTCAATCAAACCTGTATTGAAGCAATTTTCTTGCCACCAGCGACACTGCCGTTCCCCAGTCGATGCCTAACTGGT
 TGGCAGTAAGTTTGGACATAACTTCGCTTAAAGAACGGTGGGTGCGTGTGACGGCAAGGGGTGAGCTACGGATTGACCA 1600
 1601 TAATATCTCCAGCATTAACTCGATAATTTTCCCGAAATCTCTATCATCTGCTGGCGTTGATCTAATCTGTGTAGT
 ATTATAGAGGTGTAATTTAGCTATTAAAGTGGCTTTAGAGATAGTAGACGACCGCAACTAGATTAAAGACACTACTCA 1680
 1681 TTCCGATACNNNGCCATATTGGNNNCATAATCACGAACGATCACTGCATTCTGGCGTNGGGTCGGCAGCAAACATNGGCA
 AAGGCTATGNNNCGGTATAACNNNGTATTAGTGCTTGCTAGTGACGTAAGACCGCANCCAGCCGTCGTTTGTANCCGT 1760
 1761 ATGCGTGTGTAGCGGGTGAACCATTTGTCNTCGATGACGTCGGGACGCTGGTTTTACTCATCTCACGCAATACACTAACC
 TACGGACACATCGCCCACTTGGTAACAAGNAGCTACTGCAGCCCTGCGACCAAAATGAGTAGAGTGCGTTATGTGATTGC 1840
 1841 ACCCCTGGNNAACCAAGCAGCGACTGATCGCGATATTGGTACTGGGTATCCATCGCAGTGGCATACTTAAGCGTGTATATA
 TGGGGACNNNTTGGTGCTGCTGACTAGCGCTATAAGCATGACCCATAGGTAGCGTCACCGTATGAATTCGCACATATAT 1920
 1921 CTTACACTCACCGCACTGTCTTTTCGTTTGATTAAAGCATTTCCAGCACTGAAGCTAATTGACTAATACGAGTCAGGCA
 GAATGTGAGTGGCGTGACAGAAAGCAAATAATTCGCTAATAGGTGCTGACTTCGATTAAGTATTATGCTCAGTCCGT 2000

Tn insertion P7G2

2001 GCTGGGAACACCGCTCACCTCCACAGCTTTGGTACCGGTAATTTCTTTAACTCGCATCCCGGTGATGAAAGGATATTCT
 CGACCCTTGTGGCGAGTGGAGGTGTCGAAACCATGGCCATTAAAGAAATGGAGCGTAGGGCCACTACTTTCCTATAAGA 2080
 2081 GGCTGCGTAAGTAATGAATGAACCGTCCAGTAGATAAAATATTGAAAGTGATAACCTGATGTTTTAATAACGATGCAGGA
 CCGACGCATTCTACTTACTTGGCAGGTGATCTATTTATAACTTTCACTATTGGACTACAAAATTATTGCTACGTCTCT 2160
 2161 TATACATATAACATGCTGCCATCAAACCAAGTAAGCAAAATCATATTGTGCTGCCAGGTATTCAAATATCGACCGGTGG
 ATATGTATATTGTACGACGGTAGTTTGGTCCATTGCTTAGTATAACACGACGGTCCAATAAGTTTTATAGCTGCCACC 2240
 2241 TCCAGGCGGAATTTTCCACTAAATGTAGCTGTTATCAATGGGCTAATAGTAATAGCCGTATCATAGTTCTCTGAGAGCA
 AGGTCCGCCTAAAAAGGTGATTACATCGACAATAGTTACCCGATTATCATTATCGGCATAGTATCAAGAGACTCTCGT 2320
 2321 GATGTNAAAACCTCTGCTAATGGCATTGTCTGGCATAAAGGGTGAAGTCATTACCTTCCATGATAACTCATCACTCTT
 CTACANTTTTGGAGACGATTACCGTAAACAGACCGTATTTCCCACTTCAGTAATGGAAAGGTACTATTGAGTAGTGAGAA 2400
 2401 TGCTGTATTGAGTATAAATAGTAAATTAAGATTAAACGTTTATTACTACCATTTTATACCCACCCGAATAAAGTTTA
 ACGACATAACTCATATTATCATTTTAATTCTAATTGCAAAATAATGATGGTAAATATGGGGTGGGCTTATTCAAAT 2480
 2481 TGGTGATTGCGTATTACATTTTTTNAAAATGCAAGTTAAAGCCAGGTGTTTTCTATCTCAATAGCAATAAGCTCAGAGC
 ACCACTAACGCATAATGTAATAAATTTTACGTTCAATTTCGGTCCACAAAAGATAGAGTTATCGTTATTCGAGTCTCG 2560
 2561 TACTACTTGTGGTATAATAACCGTTTAAACCATCCCCATCCGCTGTGAGCTGTATAGCATAATCATGGACGTCCGGGTGT
 ATGATGAACACCATATTATTGGCAAATTGGTAGGGGTAGGCGACACTCGACATATCGTATTAGTACCTGCAGGCCACCA 2640

Tn insertion P1B9

2641 GCGCAARCRGTAGTGTCAAMTAGGCAAGCAAGGCTTAGGTAAGCTTTCCAGGTCATTTAAGAACAAAGAAATAGAAAA
 CGCGTGYGCATCACAGTKKATCCGTTCTGTTCCGAATCCATTGAAAGGTCCAGTAAATCTTGTTCTTTATCTTTTA 2720
 2721 GCTTCTGAGAAAATTTCTYCYBNN
 CGAAGACTCTTTTAAAGARGVDNN 2800

Figure 12

Sheet 2 of 5

37/39

2801 NYYKSSSCYSHKATHYYSHRWMTTAAATGGAATGCCTTTTAAACTGCCAGCATGAATCCCTCCTCAGACATAAATGGGAG 2880
 NRRMSSSGRSHMTAKRRSMYHWAATTACCTTACGGAAATTTTGACGGTCGTACTTAGGGAGGAGTCTGTATTTACCCTC
 TTTCTATCAAATTCGCTCACACCACATCCGTAAAAGCCTGATTCACATTTATTTGACTATACTCTTCTGTACAATA
 2881 AAAGATAGTTTAAGCGAGTGTGGTGTAGGCATTTTTCGACTAAGTGTAATAAAGCTGATATGAGAAGAACATGTTAT 2960
 TCAGGATGCTGTCTACATATACCTTGTACAGGCGATTCTATCATTGGGATTTTCGATAAATNNMCAATTACATTTTC
 2961 AGTCCTACGACAGATGTATATGGAACAGTGTCCGCTAAGTAGTAAGCCTAAAAGGCTATTTAANKGTTAATGTAAAAG 3040
 AGCATTGACATAAAACTTACAATTTGNAATAATTTATTAAATAAACTGTTACGATGTTTTACATCGCCATCTTATT
 3041 TCGTAAGTGTATTTTGAATGTTAAACNTTTTAAATAAATTTATTTGACAACTGCTACAAAAATGTAGCGGTAGAATAA 3120
 AAAAGTAATGTAGTCATCGACTNGGTTATATATGAAGAATTTATCTTCTAATGATAACACCATCGATTAACTMWT
 3121 TTTTTCATTAACTCAGTAGCTGACCAATATATCTTCTTAAATAGAAGGATTACTATTGTGGTAGCTAATTAGWGA 3200
 GATGAACTATATGTACTGCGATAGTGATCAAGTGCCAAAGATTTTGAACAGGCACTGGAGGGAAGCATTATGAATTT
 3201 CTACTTTGATATACATGACGCTATCAGTAGTTCACGGTTCTAAAACGTTGTCCGTTGACCTCCCTTCGTAATACTTAAA 3280
 SSTCAATCTCAAGAATACSSYSYRNHNNNNNTCTTTAGTAATCAGGCTAACTTTTTATTTTATTAACAACAATAATTWT
 3281 SSAGTTAGAGTCTTATGSSRSRYNNNNNNAGAAATCATTAGTCCGATTGAAAAATAAAAAATAATTGTTGTTATTAWA 3360
 TTGGCTGCTATCTGTGCTTACCGCAGCTTATATATCAATGGTTCRGAAACGGCAGCATATAATAGAGGATTTATCCGTTT
 3361 AACCGACGATAGACCGAATGGCGTGAATATATAGTTACCAAGYCTTTCGCGTCGTATATTATCTCTAAATAGGCAAG 3440
 TATCCGAGATGAATATTGTACTAAGCAATCAACGGTTTGAAGAAGCTGAACGTGACGCTAAAAATTTAATGTATCAATGC
 3441 ATAGGCTCTACTTATAACATGATTGTTAGTTCGCAAACTTCTTCGACTTGCACTGCGATTTTAAATTACATAGTTACG 3520
 TCATTAGCGACTGAGATTCATCATAACGATATTTTCCCTGAGGTGAGCCGGCATCTATCTGTCGGTCTTCAAATTCAC
 3521 AGTAATCGCTGACTCTAAGTAGTATTGCTATAAAAGGGACTCCACTCGGCCGTAGATAGACAGCCAGGAAGTTAACGTG 3600
 MGCCGACGCTNAACGGAGAGAAAGCACCGTCTCTTTCTGCGAGTCTCTGATATCGATGAAAAATAGCTTTCGTCGCGATAGT
 3601 KCGGCTGCGANTTGCCTCTCTTCGTGGCAGAGAAAGACGTCAGGAGACTATAGCTACTTTATCGAAAGCAGCGCTATCA 3680
 Tn insertion P3F4
 TTTATTCTTAATCATAAAAATGAGATTTTCGTTATTATCTACTGATAACCTTCAGATTATTCAACTCTACAGCCTTAAC
 3681 AAATAAGAATTAGTATTTTACTCTAAAGCAATAATAGATGACTATTGGGAAGTCTAATAAGTTGAGATGTCGGAAATTG 3760
 GCGAAAAAGCTTTCCTTTATACCCAACCATGCCGGTTTACTGGAGTGAACCAAGATACATAAACGGCAAAGGATGGC
 3761 CGCTTTTTCGAAAGGAAATATGGGTTGGGTACGGCCCAAATGACCTCACTTGGTCTTATGTATTGCCGTTTCTACCG 3840
 AACGCTTCGGTTGCGGTTGCCGATCAGGCAAGGCGTATTTTGGAGTGACGGTTAACTTCCCGATCTCATTACTAAGA
 3841 TTGCGAAGGCAACGCCAACGGCTAGTCCGTTCCGCATAAAAACTCCACTGCCAATTTGAAGGGCTAGAGTAATGATTCT 3920
 GCCACCTGCCATTAGATGATAGTATTCGAGTATGGCTGGATCAAAACAACCACTTATTGCCGTTTTCATACATCCCGGCA
 3921 CGGTGGACGGTAATCTACTATCATAAGCTCATACCGACCTAGTTTGTGGTGAATAACGGCAAAGTATGTAGGCGGT 4000
 AAAAATACGTACACAGTTAGAAAATGAACGCTGCATGATGGATGGCAGCAAATCCCGGATTTCTGATATTACGCACAA
 4001 TTTTATGCATGTGTCAATCTTTACATTGGGACGTAACCTACCGTCGTTAAGGGCCTAAAGACTATAATGCGTGT 4080
 CCTTGCATGCCCCGGATGGAGTCTGGTTACGCTGTACCCATACGGTAATCTACATAATCGCATCTAAAAATTATCCTT
 4081 GGAACGTACCGGGGCTACCTCAGACCAATGCGACATGGGTATGCCATTAGATGTATTAGCGTAGAATTTTAAATAGGAA 4160
 CAACAAATCCCTTTACATTAAACAGCATTGGTGTGATGAGTGGGCTTTTGGTGGTTACTACATCGCTCACTGGCCAA
 4161 GTTGTATTAGGGGAAATGTAATTGTGTAACCACTACTGCAGCCGAAAAACGACCAATGATGTAGCGAGTGACCGGT 4240

Figure 12

Sh et 3 of 5

38/39

ACCGTTATGGCGTTTGTGATGTCATTAATAAAACCCGAAGTGCACCGCTGAGCACACGTTTACCAGCACACGACTGG
4241 4320
TGGCAATACCGCAAAACAGCTACAGTAATTATTTTGGCGTTGACGTGGCGACTCGTGTGCAATGGTCGTGTTGCTGACC
ATGAATTAGATAGTATTGCCGGTGCTTTAAACCAACTGCTTGATCTCTACAAGTCCAATACGACAATCTGGAAAAACAA
4321 4400
TACTTAATCTATCATAACGGCCACGAAATGGTTGACGAAGTATGAGATGTTGAGGTATGCTGTTAGACCTTTTGTTT
GTCGCAGACGCACCCAGGCGCTAAATGAAGCAAAAAACGCGCTGAGCNAGCTAACAAACGTAAAGCATTCTTACG
4401 4480
CAGCGTCTGCGTGGGTCGCGATTACTTCGTTTTTTTGGCGGACTCGNTCGATTGTTGCAATTTTCGTAAGTAGAATGC
GTAATAAGTCATGAGTTACGTACTCCGATGAATGGCGTACTCGGTGCAATTGAATTATTACAAACACCCCTTTAAACAT
4481 4560
CATTATTGAGTACTCAATGCATGAGGCTACTTACCGCATGAGCCACGTTAACTTAATAATGTTGGTGGGAAATTTGTA
AGAGCAACAAGGATTAGCTGATACCGCCAGAAATGTACACTGTCTTTGTTAGCTATTATTAATAATCTGCTGGATTTTT
4561 4640
TCTCGTTGTTCTAATCGACTATGGCGGCTTTAACATGTGACAGAAACATCGATAATAATTATTAGACGACCTAAAAA
CACGCATCGAGTCTGGTCATTTACATTACATATGGAAGAAACAGCGTTACTGCCGTTACTGGACCAGGCAATGCAAAAC
4641 4720
GTGCGTAGCTCAGACCAGTAAAGTGAATGTATACCTTCTTTGTGCAATGACGGCAATGACCTGGTCCGTTACGTTTGG
ATCCAGGGGCCAGCGCNAAAGCAAAAACTGTATTACGTACTTTTGTGCGTCAACATGTCCCTCTCTATTTTCATACCG
4721 4800
TAGGTCGCCGCTCGCGNTTTCGTTTTTGTACAGTAATGCATGAAAACAGCCAGTTGTACAGGGAGAGATAAAAGTATGGC
ACAGTATCCGTTTACNNCAAATTTTGTTAATTTACTCGGGAACGCGGTAAATTTACCGAAACCGGAGGATACGTCTGA
4801 4880
TGTCTAGGCAATGNNGTTTAAACCAATTAATGAGCCCTTGCGCCATTTTAAATGGCTTTGGCCTCTATGCAGACT
CGGTCAAGCGTCATGAGGAACAATTAATTTCTGTTAGCGATAGCGGTAAAGGGATTGAAATACAGCAGCAGTCTCAA
4881 4960
GCCAGTTGCGAGTACTCCTTGTTAATTATAAAGACCAATCGCTATCGCCATTTCCCTAACCTTATGTCTGTCAGAGTT
ATCTTTACTGCTTTTATCAAGCAGACACAAATTCGCAAGGTACAGGAATTGGACTGACTATTGCGTCAAGCCTGGCTAA
4961 5040
TAGAAATGACGAAAAATAGTTCGTCTGTGTTTAAAGCGTTCATGTCTTAACTGACTGATAACGCAGTTCGGACCGATT
AATGATGGGCGGTAATCTGACACTAAAAAGTGTCGCCGGGTTGGAACCTGTGTCTCGCTAGTATTACCTTACAAGAAT
5041 5120
TTACTACCGCCATTAGACTGTGATTTTACAGGGGGCCCAACCTTGGACACAGAGCGATCATAATGGGAATGTTCTTA
In insertion
ACCAGCCGCTCAACCAATTAAGGGACGCTGTCAGNNCCGTTCTGCCTGCATCGGCAACTGGCTTGCTGGGGAATACG
5121 5200
TGGTCGGCGGAGTTGGTTAATTTCCCTGCGACAGTCNNNGGCAAGACGGACGTAGCCGTTGACCGAACGACCCCTTATGC
CGGTGAACCAACCCACCAGCAAAATGCGCTTCTCAANNCHAGAGCTTTGTATTTCTCGGAAAACTCTACGACCTGGCG
5201 5280
GCCACTTGGTGGGTGGTCGTTTTACGCGAAGAGTTNNGTCTCGAAACATAAAGAGGCCTTTGAGATGCTGGACCGC
CAACAGTTAATATTGTGTACACCAATATGCCAGTAATAAATAATTTGTTACCACCTGGCAGTTGCAGATTCTTTTGGT
5281 5360
GTTGTCAATTATAACACATGTGGTTTATACGGTCATTATTTATTAACAATGGTGGGACCGTCAACGTCTAAGAAAAACCA
TGATGATGCCGATATTAATCGGGATATCATCGGCAAAATGCTTGTCAGCCTGGGCCAACACGTCATATTGCCGCCAGTA
5361 5440
ACTACTACGGCTATAATTAGCCCTATAGTAGCCGTTTTACGAACAGTCGGACCCGTTGTGAGTGATAACGGCGGTCTAT
GTAACGAGGCTCTGACTTTATCACAACAGCAGCGATTGATTTAGTACTGATTGACATTAGAAATGCCAGAAATAGATGGT
5441 5520
CATTGCTCCGAGACTGAAATAGTGTGTGCTGCTGCTAAAGCTAAATCATGACTAACTGTAATCTTACGGTCTTTATCTACCA
ATTGAATGTGTACGATTATGGCATGATGAGCCGAATAATTTAGATCCTGACTGCATGTTTGTGGCACTATCCGCTAGCGT
5521 5600
TAACTTACACATGCTAATACCGTACTACTCGGCTTATTAATCTAGGACTGACGTACAAACACCGTGATAGGCGATCGCA
ASCNVHAGAWRHTMTCTRTYGTDDAAAAAAHRRDGRKDHTCATHAYANHTTACAAAACCAAGTGACATTGGCTACCTTAGC
5601 5680
TSGBNKCTWYAKWAGYARCAHHTTTTTTHYHCYHMDWAGTADTRTNNAATGTTTGGTCACTGTAACCGATGGAATCG

Figure 12

Sh et 4 of 5

39/39

5681 TCGCTACATCAGTATTGCCGCAGAATACCAACTTTACGAAATATAGAGCTACAGGAGCAGGATCC 5746

AGCGATGTAGTCATAACGGCGTCTTATGGTTGAAAATGCTTTATATCTCGATGTCTCGTCCTAGG